



Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization

Fabian Pedregosa, Rémi Leblond, Simon Lacoste-Julien

► To cite this version:

Fabian Pedregosa, Rémi Leblond, Simon Lacoste-Julien. Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization. NIPS 2017 - Thirty-First Annual Conference on Neural Information Processing Systems, Dec 2017, Long Beach, United States. pp.1-28. hal-01638058

HAL Id: hal-01638058

<https://inria.hal.science/hal-01638058>

Submitted on 19 Nov 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization

Fabian Pedregosa
INRIA/ENS*
Paris, France

Rémi Leblond
INRIA/ENS*
Paris, France

Simon Lacoste-Julien
MILA and DIRO
Université de Montréal, Canada

Abstract

Due to their simplicity and excellent performance, parallel asynchronous variants of stochastic gradient descent have become popular methods to solve a wide range of large-scale optimization problems on multi-core architectures. Yet, despite their practical success, support for nonsmooth objectives is still lacking, making them unsuitable for many problems of interest in machine learning, such as the Lasso, group Lasso or empirical risk minimization with convex constraints. In this work, we propose and analyze PROXASAGA, a fully asynchronous sparse method inspired by SAGA, a variance reduced incremental gradient algorithm. The proposed method is easy to implement and significantly outperforms the state of the art on several nonsmooth, large-scale problems. We prove that our method achieves a theoretical linear speedup with respect to the sequential version under assumptions on the sparsity of gradients and block-separability of the proximal term. Empirical benchmarks on a multi-core architecture illustrate practical speedups of up to 12x on a 20-core machine.

1 Introduction

The widespread availability of multi-core computers motivates the development of parallel methods adapted for these architectures. One of the most popular approaches is HOGWILD (Niu et al., 2011), an asynchronous variant of stochastic gradient descent (SGD). In this algorithm, multiple threads run the update rule of SGD asynchronously in parallel. As SGD, it only requires visiting a small batch of random examples per iteration, which makes it ideally suited for large scale machine learning problems. Due to its simplicity and excellent performance, this parallelization approach has recently been extended to other variants of SGD with better convergence properties, such as SVRG (Johnson & Zhang, 2013) and SAGA (Defazio et al., 2014).

Despite their practical success, existing parallel asynchronous variants of SGD are limited to smooth objectives, making them inapplicable to many problems in machine learning and signal processing. In this work, we develop a sparse variant of the SAGA algorithm and consider its parallel asynchronous variants for general *composite* optimization problems of the form:

$$\arg \min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) + h(\mathbf{x}) \quad , \quad \text{with } f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \quad , \quad (\text{OPT})$$

where each f_i is convex with L -Lipschitz gradient, the average function f is μ -strongly convex and h is convex but potentially nonsmooth. We further assume that h is “simple” in the sense that we have access to its proximal operator, and that it is block-separable, that is, it can be decomposed block coordinate-wise as $h(\mathbf{x}) = \sum_{B \in \mathcal{B}} h_B([\mathbf{x}]_B)$, where \mathcal{B} is a partition of the coefficients into

*DI École normale supérieure, CNRS, PSL Research University

subsets which will call *blocks* and h_B only depends on coordinates in block B . Note that there is no loss of generality in this last assumption as a unique block covering all coordinates is a valid partition, though in this case, our sparse variant of the SAGA algorithm reduces to the original SAGA algorithm and no gain from sparsity is obtained.

This template models a broad range of problems arising in machine learning and signal processing: the finite-sum structure of f includes the least squares or logistic loss functions; the proximal term h includes penalties such as the ℓ_1 or group lasso penalty. Furthermore, this term can be extended-valued, thus allowing for convex constraints through the indicator function.

Contributions. This work presents two main contributions. First, in §2 we describe Sparse Proximal SAGA, a novel variant of the SAGA algorithm which features a reduced cost per iteration in the presence of sparse gradients and a block-separable penalty. Like other variance reduced methods, it enjoys a linear convergence rate under strong convexity. Second, in §3 we present PROXASAGA, a lock-free asynchronous parallel version of the aforementioned algorithm that does not require consistent reads. Our main results states that PROXASAGA obtains (under assumptions) a theoretical linear speedup with respect to its sequential version. Empirical benchmarks reported in §4 show that this method dramatically outperforms state-of-the-art alternatives on large sparse datasets, while the empirical speedup analysis illustrates the practical gains as well as its limitations.

1.1 Related work

Asynchronous coordinate-descent. For composite objective functions of the form (OPT), most of the existing literature on asynchronous optimization has focused on variants of coordinate descent. Liu & Wright (2015) proposed an asynchronous variant of (proximal) coordinate descent and proved a near-linear speedup in the number of cores used, given a suitable step size. This approach has been recently extended to general block-coordinate schemes by Peng et al. (2016), to greedy coordinate-descent schemes by You et al. (2016) and to non-convex problems by Davis et al. (2016). However, as illustrated by our experiments, in the large sample regime coordinate descent compares poorly against incremental gradient methods like SAGA.

Variance reduced incremental gradient and their asynchronous variants. Initially proposed in the context of smooth optimization by Le Roux et al. (2012), variance reduced incremental gradient methods have since been extended to minimize composite problems of the form (OPT) (see table below). Smooth variants of these methods have also recently been extended to the asynchronous setting, where multiple threads run the update rule asynchronously and in parallel. Interestingly, none of these methods achieve both simultaneously, i.e. asynchronous optimization of composite problems. Since variance reduced incremental gradient methods have shown state of the art performance in both settings, this generalization is of key practical interest.

Objective	Sequential Algorithm	Asynchronous Algorithm
Smooth	SVRG (Johnson & Zhang, 2013)	SVRG (Reddi et al., 2015)
	SDCA (Shalev-Shwartz & Zhang, 2013)	PASSCODE (Hsieh et al., 2015, SDCA variant)
	SAGA (Defazio et al., 2014)	ASAGA (Leblond et al., 2017, SAGA variant)
Composite	PROXSDCA (Shalev-Shwartz et al., 2012)	
	SAGA (Defazio et al., 2014)	This work: PROXASAGA
	ProxSVRG (Xiao & Zhang, 2014)	

On the difficulty of a composite extension. Two key issues explain the paucity in the development of asynchronous incremental gradient methods for composite optimization. The first issue is related to the design of such algorithms. Asynchronous variants of SGD are most competitive when the updates are sparse and have a small overlap, that is, when each update modifies a small and different subset of the coefficients. This is typically achieved by updating only coefficients for which the partial gradient at a given iteration is nonzero,² but existing schemes such as the lagged updates technique (Schmidt et al., 2016) are not applicable in the asynchronous setting. The second

²Although some regularizers are sparsity inducing, large scale datasets are often extremely sparse and leveraging this property is crucial for the efficiency of the method.

difficulty is related to the analysis of such algorithms. All convergence proofs crucially use the Lipschitz condition on the gradient to bound the noise terms derived from asynchrony. However, in the composite case, the gradient mapping term (Beck & Teboulle, 2009), which replaces the gradient in proximal-gradient methods, does not have a bounded Lipschitz constant. Hence, the traditional proof technique breaks down in this scenario.

Other approaches. Recently, Meng et al. (2017); Gu et al. (2016) independently proposed a doubly stochastic method to solve the problem at hand. Following Meng et al. (2017) we refer to it as Async-PROXSVRCD. This method performs coordinate descent-like updates in which the true gradient is replaced by its SVRG approximation. It hence features a doubly-stochastic loop: at each iteration we select a random coordinate *and* a random sample. Because the selected coordinate block is uncorrelated with the chosen sample, the algorithm can be orders of magnitude slower than SAGA in the presence of sparse gradients. Appendix F contains a comparison of these methods.

1.2 Definitions and notations

By convention, we denote vectors and vector-valued functions in lowercase boldface (e.g. \mathbf{x}) and matrices in uppercase boldface (e.g. \mathbf{D}). The proximal operator of a convex lower semicontinuous function h is defined as $\text{prox}_h(\mathbf{x}) := \arg \min_{\mathbf{z} \in \mathbb{R}^p} \{h(\mathbf{z}) + \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|^2\}$. A function f is said to be L -smooth if it is differentiable and its gradient is L -Lipschitz continuous. A function f is said to be μ -strongly convex if $f - \frac{\mu}{2}\|\cdot\|^2$ is convex. We use the notation $\kappa := L/\mu$ to denote the condition number for an L -smooth and μ -strongly convex function.³

\mathbf{I}_p denotes the p -dimensional identity matrix, $\mathbb{1}\{\text{cond}\}$ the characteristic function, which is 1 if cond evaluates to true and 0 otherwise. The average of a vector or matrix is denoted $\bar{\alpha} := \frac{1}{n} \sum_{i=1}^n \alpha_i$. We use $\|\cdot\|$ for the Euclidean norm. For a positive semi-definite matrix \mathbf{D} , we define its associated distance as $\|\mathbf{x}\|_{\mathbf{D}}^2 := \langle \mathbf{x}, \mathbf{D}\mathbf{x} \rangle$. We denote by $[\mathbf{x}]_b$ the b -th coordinate in \mathbf{x} . This notation is overloaded so that for a collection of blocks $T = \{B_1, B_2, \dots\}$, $[\mathbf{x}]_T$ denotes the vector \mathbf{x} restricted to the coordinates in the blocks of T . For convenience, when T consists of a single block B we use $[\mathbf{x}]_B$ as a shortcut of $[\mathbf{x}]_{\{B\}}$. Finally, we distinguish \mathbb{E} , the full expectation taken with respect to all the randomness in the system, from \mathbb{E} , the conditional expectation of a random i_t (the random feature sampled at each iteration by SGD-like algorithms) conditioned on all the “past”, which the context will clarify.

2 Sparse Proximal SAGA

Original SAGA algorithm. The original SAGA algorithm (Defazio et al., 2014) maintains two moving quantities: the current iterate \mathbf{x} and a table (memory) of historical gradients $(\alpha_i)_{i=1}^n$. At every iteration, it samples an index $i \in \{1, \dots, n\}$ uniformly at random, and computes the next iterate (\mathbf{x}^+, α^+) according to the following recursion:

$$\mathbf{u}_i = \nabla f_i(\mathbf{x}) - \alpha_i + \bar{\alpha}; \quad \mathbf{x}^+ = \text{prox}_{\gamma h}(\mathbf{x} - \gamma \mathbf{u}_i); \quad \alpha_i^+ = \nabla f_i(\mathbf{x}). \quad (1)$$

On each iteration, this update rule requires to visit all coefficients even if the partial gradients ∇f_i are sparse. Sparse partial gradients arise in a variety of practical scenarios: for example, in generalized linear models the partial gradients inherit the sparsity pattern of the dataset. Given that large-scale datasets are often sparse,⁴ leveraging this sparsity is crucial for the success of the optimizer.

Sparse Proximal SAGA algorithm. We will now describe an algorithm that leverages sparsity in the partial gradients by only updating those blocks that intersect with the support of the partial gradients. Since in this update scheme some blocks might appear more frequently than others, we will need to counterbalance this undesirable effect with a well-chosen block-wise reweighting of the average gradient and the proximal term.

In order to make precise this block-wise reweighting, we define the following quantities. We denote by T_i the *extended support* of ∇f_i , which is the set of blocks that intersect the support of ∇f_i ,

³Since we have assumed that each individual f_i is L -smooth, f itself is L -smooth – but it could have a smaller smoothness constant. Our rates are in terms of this bigger L/μ , as is standard in the SAGA literature.

⁴For example, in the LibSVM datasets suite, 8 out of the 11 datasets (as of May 2017) with more than a million samples have a density between 10^{-4} and 10^{-6} .

formally defined as $T_i := \{B : \text{supp}(\nabla f_i) \cap B \neq \emptyset, B \in \mathcal{B}\}$. For totally separable penalties such as the ℓ_1 norm, the blocks are individual coordinates and so the extended support covers the same coordinates as the support. Let $d_B := n/n_B$, where $n_B := \sum_i \mathbb{1}\{B \in T_i\}$ is the number of times that $B \in T_i$. For simplicity we assume $n_B > 0$, as otherwise the problem can be reformulated without block B . The update rule in (1) requires computing the proximal operator of h , which involves a full pass on the coordinates. In our proposed algorithm, we replace h in (1) with the function $\varphi_i(\mathbf{x}) := \sum_{B \in T_i} d_B h_B(\mathbf{x})$, whose form is justified by the following three properties. First, this function is zero outside T_i , allowing for sparse updates. Second, because of the block-wise reweighting d_B , the function φ_i is an unbiased estimator of h (i.e., $\mathbf{E} \varphi_i = h$), property which will be crucial to prove the convergence of the method. Third, φ_i inherits the block-wise structure of h and its proximal operator can be computed from that of h as $[\text{prox}_{\gamma \varphi_i}(\mathbf{x})]_B = [\text{prox}_{(d_B \gamma) h_B}(\mathbf{x})]_B$ if $B \in T_i$ and $[\text{prox}_{\gamma \varphi_i}(\mathbf{x})]_B = [\mathbf{x}]_B$ otherwise. Following Leblond et al. (2017), we will also replace the dense gradient estimate \mathbf{u}_i by the sparse estimate $\mathbf{v}_i := \nabla f_i(\mathbf{x}) - \boldsymbol{\alpha}_i + \mathbf{D}_i \bar{\boldsymbol{\alpha}}$, where \mathbf{D}_i is the diagonal matrix defined block-wise as $[\mathbf{D}_i]_{B,B} = d_B \mathbb{1}\{B \in T_i\} \mathbf{I}_{|B|}$. It is easy to verify that the vector $\mathbf{D}_i \bar{\boldsymbol{\alpha}}$ is a weighted projection onto the support of T_i and $\mathbf{E} \mathbf{D}_i \bar{\boldsymbol{\alpha}} = \bar{\boldsymbol{\alpha}}$, making \mathbf{v}_i an unbiased estimate of the gradient.

We now have all necessary elements to describe the Sparse Proximal SAGA algorithm. As the original SAGA algorithm, it maintains two moving quantities: the current iterate $\mathbf{x} \in \mathbb{R}^p$ and a table of historical gradients $(\boldsymbol{\alpha}_i)_{i=1}^n$, $\boldsymbol{\alpha}_i \in \mathbb{R}^p$. At each iteration, the algorithm samples an index $i \in \{1, \dots, n\}$ and computes the next iterate $(\mathbf{x}^+, \boldsymbol{\alpha}^+)$ as:

$$\mathbf{v}_i = \nabla f_i(\mathbf{x}) - \boldsymbol{\alpha}_i + \mathbf{D}_i \bar{\boldsymbol{\alpha}}; \mathbf{x}^+ = \text{prox}_{\gamma \varphi_i}(\mathbf{x} - \gamma \mathbf{v}_i); \boldsymbol{\alpha}_i^+ = \nabla f_i(\mathbf{x}), \quad (\text{SPS})$$

where in a practical implementation the vector $\bar{\boldsymbol{\alpha}}$ is updated incrementally at each iteration.

The above algorithm is sparse in the sense that it only requires to visit and update blocks in the extended support: if $B \notin T_i$, by the sparsity of \mathbf{v}_i and prox_{φ_i} , we have $[\mathbf{x}^+]_B = [\mathbf{x}]_B$. Hence, when the extended support T_i is sparse, this algorithm can be orders of magnitude faster than the naive SAGA algorithm. The extended support is sparse for example when the partial gradients are sparse and the penalty is separable, as is the case of the ℓ_1 norm or the indicator function over a hypercube, or when the the penalty is block-separable in a way such that only a small subset of the blocks overlap with the support of the partial gradients. Initialization of variables and a reduced storage scheme for the memory are discussed in the implementation details section of Appendix E.

Relationship with existing methods. This algorithm can be seen as a generalization of both the Standard SAGA algorithm and the Sparse SAGA algorithm of Leblond et al. (2017). When the proximal term is not block-separable, then $d_B = 1$ (for a unique block B) and the algorithm defaults to the Standard (dense) SAGA algorithm. In the smooth case (i.e., $h = 0$), the algorithm defaults to the Sparse SAGA method. Hence we note that the sparse gradient estimate \mathbf{v}_i in our algorithm is the same as the one proposed in Leblond et al. (2017). However, we emphasize that a straightforward combination of this sparse update rule with the proximal update from the Standard SAGA algorithm results in a nonconvergent algorithm: the block-wise reweighting of h is a surprisingly simple but crucial change. We now give the convergence guarantees for this algorithm.

Theorem 1. *Let $\gamma = \frac{a}{5L}$ for any $a \leq 1$ and f be μ -strongly convex ($\mu > 0$). Then Sparse Proximal SAGA converges geometrically in expectation with a rate factor of at least $\rho = \frac{1}{5} \min\{\frac{1}{n}, \frac{1}{\kappa}\}$. That is, for \mathbf{x}_t obtained after t updates, we have the following bound:*

$$\mathbb{E} \|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq (1 - \rho)^t C_0, \text{ with } C_0 := \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{1}{5L^2} \sum_{i=1}^n \|\boldsymbol{\alpha}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2.$$

Remark. For the step size $\gamma = 1/5L$, the convergence rate is $(1 - 1/5 \min\{1/n, 1/\kappa\})$. We can thus identify two regimes: the “big data” regime, $n \geq \kappa$, in which the rate factor is bounded by $1/5n$, and the “ill-conditioned” regime, $\kappa \geq n$, in which the rate factor is bounded by $1/5\kappa$. This rate roughly matches the rate obtained by Defazio et al. (2014). While the step size bound of $1/5L$ is slightly smaller than the $1/3L$ one obtained in that work, this can be explained by their stronger assumptions: each f_i is strongly convex whereas they are strongly convex only on average in this work. All proofs for this section can be found in Appendix B.

Algorithm 1 PROXASAGA (analyzed)	Algorithm 2 PROXASAGA (implemented)
1: Initialize shared variables \mathbf{x} and $(\alpha_i)_{i=1}^n$ 2: keep doing in parallel 3: $\hat{\mathbf{x}} =$ inconsistent read of \mathbf{x} 4: $\hat{\alpha} =$ inconsistent read of α 5: Sample i uniformly in $\{1, \dots, n\}$ 6: $S_i :=$ support of ∇f_i 7: $T_i :=$ extended support of ∇f_i in \mathcal{B} 8: $[\bar{\alpha}]_{T_i} = 1/n \sum_{j=1}^n [\hat{\alpha}_j]_{T_i}$ 9: $[\delta\alpha]_{S_i} = [\nabla f_i(\hat{\mathbf{x}})]_{S_i} - [\hat{\alpha}_i]_{S_i}$ 10: $[\hat{\mathbf{v}}]_{T_i} = [\delta\alpha]_{T_i} + [\mathbf{D}_i \bar{\alpha}]_{T_i}$ 11: $[\delta\mathbf{x}]_{T_i} = [\text{prox}_{\gamma\varphi_i}(\hat{\mathbf{x}} - \gamma\hat{\mathbf{v}})]_{T_i} - [\hat{\mathbf{x}}]_{T_i}$ 12: for B in T_i do 13: for $b \in B$ do 14: $[\mathbf{x}]_b \leftarrow [\mathbf{x}]_b + [\delta\mathbf{x}]_b$ \triangleright atomic 15: if $b \in S_i$ then 16: $[\alpha_i]_b \leftarrow [\nabla f_i(\hat{\mathbf{x}})]_b$ 17: end if 18: end for 19: end for 20: // (' \leftarrow ' denotes shared memory update.) 21: end parallel loop	1: Initialize shared variables \mathbf{x} , $(\alpha_i)_{i=1}^n$, $\bar{\alpha}$ 2: keep doing in parallel 3: Sample i uniformly in $\{1, \dots, n\}$ 4: $S_i :=$ support of ∇f_i 5: $T_i :=$ extended support of ∇f_i in \mathcal{B} 6: $[\hat{\mathbf{x}}]_{T_i} =$ inconsistent read of \mathbf{x} on T_i 7: $\hat{\alpha}_i =$ inconsistent read of α_i 8: $[\bar{\alpha}]_{T_i} =$ inconsistent read of $\bar{\alpha}$ on T_i 9: $[\delta\alpha]_{S_i} = [\nabla f_i(\hat{\mathbf{x}})]_{S_i} - [\hat{\alpha}_i]_{S_i}$ 10: $[\hat{\mathbf{v}}]_{T_i} = [\delta\alpha]_{T_i} + [\mathbf{D}_i \bar{\alpha}]_{T_i}$ 11: $[\delta\mathbf{x}]_{T_i} = [\text{prox}_{\gamma\varphi_i}(\hat{\mathbf{x}} - \gamma\hat{\mathbf{v}})]_{T_i} - [\hat{\mathbf{x}}]_{T_i}$ 12: for B in T_i do 13: for $b \in B$ do 14: $[\mathbf{x}]_b \leftarrow [\mathbf{x}]_b + [\delta\mathbf{x}]_b$ \triangleright atomic 15: if $b \in S_i$ then 16: $[\bar{\alpha}]_b \leftarrow [\bar{\alpha}]_b + 1/n[\delta\alpha]_b$ \triangleright atomic 17: end if 18: end for 19: end for 20: $\alpha_i \leftarrow \nabla f_i(\hat{\mathbf{x}})$ (scalar update) \triangleright atomic 21: end parallel loop

3 Asynchronous Sparse Proximal SAGA

We introduce PROXASAGA – the asynchronous parallel variant of Sparse Proximal SAGA. In this algorithm, multiple cores update a central parameter vector using the Sparse Proximal SAGA introduced in the previous section, and updates are performed asynchronously. The algorithm parameters are read and written without vector locks, i.e., the vector content of the shared memory can potentially change while a core is reading or writing to main memory coordinate by coordinate. These operations are typically called *inconsistent* (at the vector level).

The full algorithm is described in Algorithm 1 for its theoretical version (on which our analysis is built) and in Algorithm 2 for its practical implementation. The practical implementation differs from the analyzed algorithm in three points. First, in the implemented algorithm, index i is sampled before reading the coefficients to minimize memory access since only the extended support needs to be read. Second, since our implementation targets generalized linear models, the memory α_i can be compressed into a single scalar in L20 (see Appendix E). Third, $\bar{\alpha}$ is stored in memory and updated incrementally instead of recomputed at each iteration.

The rest of the section is structured as follows: we start by describing our framework of analysis; we then derive essential properties of PROXASAGA along with a classical delay assumption. Finally, we state our main convergence and speedup result.

3.1 Analysis framework

As in most of the recent asynchronous optimization literature, we build on the hardware model introduced by Niu et al. (2011), with multiple cores reading and writing to a shared memory parameter vector. These operations are asynchronous (lock-free) and *inconsistent*:⁵ $\hat{\mathbf{x}}_t$, the local copy of the parameters of a given core, does not necessarily correspond to a consistent iterate in memory.

“Perturbed” iterates. To handle this additional difficulty, contrary to most contributions in this field, we choose the “perturbed iterate framework” proposed by Mania et al. (2017) and refined by Leblond et al. (2017). This framework can analyze variants of SGD which obey the update rule:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \mathbf{v}(\mathbf{x}_t, i_t), \text{ where } \mathbf{v} \text{ verifies the unbiasedness condition } \mathbf{E} \mathbf{v}(\mathbf{x}, i_t) = \nabla f(\mathbf{x})$$

⁵This is an extension of the framework of Niu et al. (2011), where consistent updates were assumed.

and the expectation is computed with respect to i_t . In the asynchronous parallel setting, cores are reading inconsistent iterates from memory, which we denote \hat{x}_t . As these inconsistent iterates are affected by various delays induced by asynchrony, they cannot easily be written as a function of their previous iterates. To alleviate this issue, Mania et al. (2017) choose to introduce an additional quantity for the purpose of the analysis:

$$x_{t+1} := x_t - \gamma v(\hat{x}_t, i_t), \quad \text{the “virtual iterate” – which is never actually computed.} \quad (2)$$

Note that this equation is the *definition* of this new quantity x_t . This virtual iterate is useful for the convergence analysis and makes for much easier proofs than in the related literature.

“After read” labeling. How we choose to define the iteration counter t to label an iterate x_t matters in the analysis. In this paper, we follow the “after read” labeling proposed in Leblond et al. (2017), in which we update our iterate counter, t , as each core *finishes reading* its copy of the parameters (in the specific case of PROXASAGA, this includes both \hat{x}_t and $\hat{\alpha}^t$). This means that \hat{x}_t is the $(t + 1)^{th}$ fully completed read. One key advantage of this approach compared to the classical choice of Niu et al. (2011) – where t is increasing after each successful update – is that it guarantees both that the i_t are uniformly distributed and that i_t and \hat{x}_t are independent. This property is not verified when using the “after write” labeling of Niu et al. (2011), although it is still implicitly assumed in the papers using this approach, see Leblond et al. (2017, Section 3.2) for a discussion of issues related to the different labeling schemes.

Generalization to composite optimization. Although the perturbed iterate framework was designed for gradient-based updates, we can extend it to proximal methods by remarking that in the sequential setting, proximal stochastic gradient descent and its variants can be characterized by the following similar update rule:

$$x_{t+1} = x_t - \gamma g(x_t, v_{i_t}, i_t), \quad \text{with } g(x, v, i) := \frac{1}{\gamma} (x - \text{prox}_{\gamma \varphi_i}(x - \gamma v)), \quad (3)$$

where as before v verifies the unbiasedness condition $\mathbb{E} v = \nabla f(x)$. The Proximal Sparse SAGA iteration can be easily written within this template by using φ_i and v_i as defined in §2. Using this definition of g , we can define PROXASAGA virtual iterates as:

$$x_{t+1} := x_t - \gamma g(\hat{x}_t, \hat{v}_{i_t}^t, i_t), \quad \text{with } \hat{v}_{i_t}^t = \nabla f_{i_t}(\hat{x}_t) - \hat{\alpha}_{i_t}^t + D_{i_t} \bar{\alpha}^t, \quad (4)$$

where as in the sequential case, the memory terms are updated as $\hat{\alpha}_{i_t}^t = \nabla f_{i_t}(\hat{x}_t)$. Our theoretical analysis of PROXASAGA will be based on this definition of the virtual iterate x_{t+1} .

3.2 Properties and assumptions

Now that we have introduced the “after read” labeling for proximal methods in Eq. (4), we can leverage the framework of Leblond et al. (2017, Section 3.3) to derive essential properties for the analysis of PROXASAGA. We describe below three useful properties arising from the definition of Algorithm 1, and then state a central (but standard) assumption that the delays induced by the asynchrony are uniformly bounded.

Independence: Due to the “after read” global ordering, i_r is independent of \hat{x}_t for all $r \geq t$. We enforce the independence for $r = t$ by having the cores read all the shared parameters before their iterations.

Unbiasedness: The term $\hat{v}_{i_t}^t$ is an unbiased estimator of the gradient of f at \hat{x}_t . This property is a consequence of the independence between i_t and \hat{x}_t .

Atomicity: The shared parameter coordinate update of $[x]_b$ on Line 14 is atomic. This means that there are no overwrites for a single coordinate even if several cores compete for the same resources. Most modern processors have support for atomic operations with minimal overhead.

Bounded overlap assumption. We assume that there exists a uniform bound, τ , on the maximum number of overlapping iterations. This means that every coordinate update from iteration t is successfully written to memory before iteration $t + \tau + 1$ starts. Our result will give us conditions on τ to obtain linear speedups.

Bounding $\hat{x}_t - x_t$. The delay assumption of the previous paragraph allows to express the difference between real and virtual iterate using the gradient mapping $g_u := g(\hat{x}_u, \hat{v}_{i_u}^u, i_u)$ as:

$$\hat{x}_t - x_t = \gamma \sum_{u=(t-\tau)_+}^{t-1} G_u^t g_u, \quad \text{where } G_u^t \text{ are } p \times p \text{ diagonal matrices with terms in } \{0, +1\}. \quad (5)$$

0 represents instances where both \hat{x}_u and x_u have received the corresponding updates. +1, on the contrary, represents instances where \hat{x}_u has not yet received an update that is already in x_u by definition. This bound will prove essential to our analysis.

3.3 Analysis

In this section, we state our convergence and speedup results for PROXASAGA. The full details of the analysis can be found in [Appendix C](#). Following [Niu et al. \(2011\)](#), we introduce a sparsity measure (generalized to the composite setting) that will appear in our results.

Definition 1. Let $\Delta := \max_{B \in \mathcal{B}} |\{i : T_i \ni B\}|/n$. This is the normalized maximum number of times that a block appears in the extended support. For example, if a block is present in all T_i , then $\Delta = 1$. If no two T_i share the same block, then $\Delta = 1/n$. We always have $1/n \leq \Delta \leq 1$.

Theorem 2 (Convergence guarantee of PROXASAGA). Suppose $\tau \leq \frac{1}{10\sqrt{\Delta}}$. For any step size $\gamma = \frac{a}{L}$ with $a \leq a^*(\tau) := \frac{1}{36} \min\{1, \frac{6\kappa}{\tau}\}$, the inconsistent read iterates of Algorithm 1 converge in expectation at a geometric rate factor of at least: $\rho(a) = \frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$, i.e. $\mathbb{E}\|\hat{x}_t - x^*\|^2 \leq (1 - \rho)^t \tilde{C}_0$, where \tilde{C}_0 is a constant independent of t ($\approx \frac{n\kappa}{a} C_0$ with C_0 as defined in Theorem 1).

This last result is similar to the original SAGA convergence result and our own Theorem 1, with both an extra condition on τ and on the maximum allowable step size. In the best sparsity case, $\Delta = 1/n$ and we get the condition $\tau \leq \sqrt{n}/10$. We now compare the geometric rate above to the one of Sparse Proximal SAGA to derive the necessary conditions under which PROXASAGA is linearly faster.

Corollary 1 (Speedup). Suppose $\tau \leq \frac{1}{10\sqrt{\Delta}}$. If $\kappa \geq n$, then using the step size $\gamma = 1/36L$, PROXASAGA converges geometrically with rate factor $\Omega(\frac{1}{\kappa})$. If $\kappa < n$, then using the step size $\gamma = 1/36n\mu$, PROXASAGA converges geometrically with rate factor $\Omega(\frac{1}{n})$. In both cases, the convergence rate is the same as Sparse Proximal SAGA. Thus PROXASAGA is linearly faster than its sequential counterpart up to a constant factor. Note that in both cases the step size does not depend on τ .

Furthermore, if $\tau \leq 6\kappa$, we can use a universal step size of $\Theta(1/L)$ to get a similar rate for PROXASAGA than Sparse Proximal SAGA, thus making it adaptive to local strong convexity since the knowledge of κ is not required.

These speedup regimes are comparable with the best ones obtained in the smooth case, including [Niu et al. \(2011\)](#); [Reddi et al. \(2015\)](#), even though unlike these papers, we support inconsistent reads and nonsmooth objective functions. The one exception is [Leblond et al. \(2017\)](#), where the authors prove that their algorithm, ASAGA, can obtain a linear speedup even without sparsity in the well-conditioned regime. In contrast, PROXASAGA always requires some sparsity. Whether this property for smooth objective functions could be extended to the composite case remains an open problem.

Relative to ASYSPCD, in the best case scenario (where the components of the gradient are uncorrelated, a somewhat unrealistic setting), ASYSPCD can get a near-linear speedup for τ as big as $\sqrt[4]{p}$. Our result states that $\tau = \mathcal{O}(1/\sqrt{\Delta})$ is necessary for a linear speedup. This means in case $\Delta \leq 1/\sqrt{p}$ our bound is better than the one obtained for ASYSPCD. Recalling that $1/n \leq \Delta \leq 1$, it appears that PROXASAGA is favored when n is bigger than \sqrt{p} whereas ASYSPCD may have a better bound otherwise, though this comparison should be taken with a grain of salt given the assumptions we had to make to arrive at comparable quantities. An extended comparison with the related work can be found in [Appendix D](#).

4 Experiments

In this section, we compare PROXASAGA with related methods on different datasets. Although PROXASAGA can be applied more broadly, we focus on $\ell_1 + \ell_2$ -regularized logistic regression, a model of particular practical importance. The objective function takes the form

$$\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i \mathbf{a}_i^\top \mathbf{x})) + \frac{\lambda_1}{2} \|\mathbf{x}\|_2^2 + \lambda_2 \|\mathbf{x}\|_1, \quad (6)$$

where $\mathbf{a}_i \in \mathbb{R}^p$ and $b_i \in \{-1, +1\}$ are the data samples. Following [Defazio et al. \(2014\)](#), we set $\lambda_1 = 1/n$. The amount of ℓ_1 regularization (λ_2) is selected to give an approximate $1/10$ nonzero

Table 1: Description of datasets.

Dataset	n	p	density	L	Δ
KDD 2010 (Yu et al., 2010)	19,264,097	1,163,024	10^{-6}	28.12	0.15
KDD 2012 (Juan et al., 2016)	149,639,105	54,686,452	2×10^{-7}	1.25	0.85
Criteo (Juan et al., 2016)	45,840,617	1,000,000	4×10^{-5}	1.25	0.89

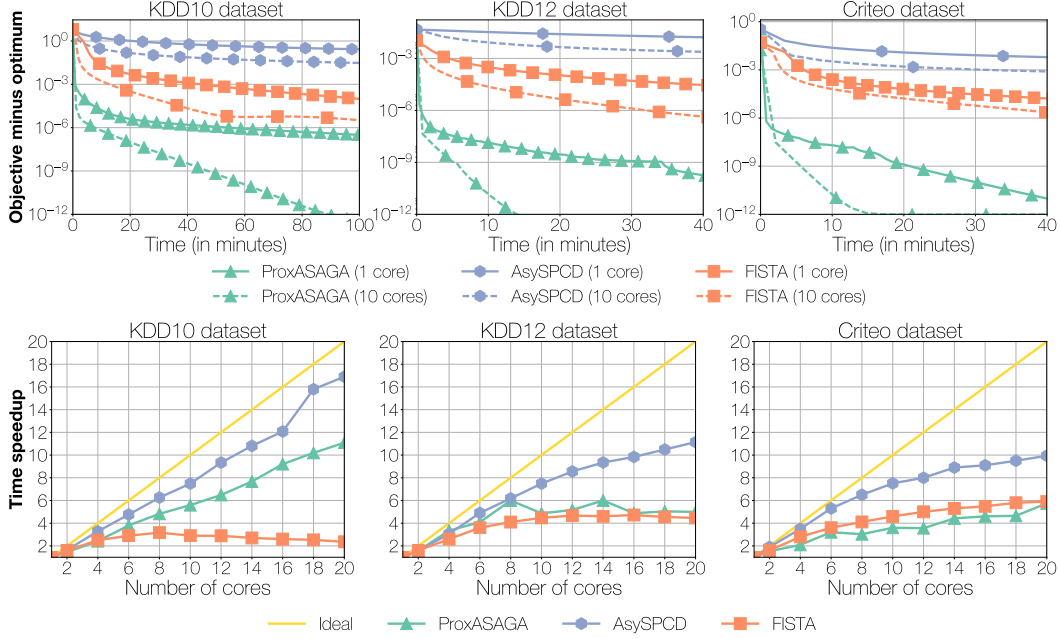


Figure 1: **Convergence for asynchronous stochastic methods for $\ell_1 + \ell_2$ -regularized logistic regression.** Top: Suboptimality as a function of time for different asynchronous methods using 1 and 10 cores. Bottom: Running time speedup as function of the number of cores. PROXASAGA achieves significant speedups over its sequential version while being orders of magnitude faster than competing methods. ASYSPCD achieves the highest speedups but it also the slowest overall method.

coefficients. Implementation details are available in [Appendix E](#). We chose the 3 datasets described in Table 1

Results. We compare three parallel asynchronous methods on the aforementioned datasets: PROXASAGA (this work),⁶ ASYSPCD, the asynchronous proximal coordinate descent method of Liu & Wright (2015) and the (synchronous) FISTA algorithm (Beck & Teboulle, 2009), in which the gradient computation is parallelized by splitting the dataset into equal batches. We aim to benchmark these methods in the most realistic scenario possible; to this end we use the following step size: $1/2L$ for PROXASAGA, $1/L_c$ for ASYSPCD, where L_c is the coordinate-wise Lipschitz constant of the gradient, while FISTA uses backtracking line-search. The results can be seen in Figure 1 (top) with both one (thus sequential) and ten processors. Two main observations can be made from this figure. First, PROXASAGA is significantly faster on these problems. Second, its asynchronous version offers a significant speedup over its sequential counterpart.

In Figure 1 (bottom) we present speedup with respect to the number of cores, where speedup is computed as the time to achieve a suboptimality of 10^{-10} with one core divided by the time to achieve the same suboptimality using several cores. While our *theoretical speedups* (with respect to the number of iterations) are almost linear as our theory predicts (see [Appendix F](#)), we observe a different story for our *running time* speedups. This can be attributed to memory access overhead, which our model does not take into account. As predicted by our theoretical results, we observe

⁶A reference C++/Python implementation of is available at <https://github.com/fabianp/ProxASAGA>

a high correlation between the Δ dataset sparsity measure and the empirical speedup: KDD 2010 ($\Delta = 0.15$) achieves a 11x speedup, while in Criteo ($\Delta = 0.89$) the speedup is never above 6x.

Note that although competitor methods exhibit similar or sometimes better speedups, they remain orders of magnitude slower than PROXASAGA in running time for large sparse problems. In fact, our method is between 5x and 80x times faster (in time to reach 10^{-10} suboptimality) than FISTA and between 13x and 290x times faster than ASYSPCD (see [Appendix F.3](#)).

5 Conclusion and future work

In this work, we have described PROXASAGA, an asynchronous variance reduced algorithm with support for composite objective functions. This method builds upon a novel sparse variant of the (proximal) SAGA algorithm that takes advantage of sparsity in the individual gradients. We have proven that this algorithm is linearly convergent under a condition on the step size and that it is linearly faster than its sequential counterpart given a bound on the delay. Empirical benchmarks show that PROXASAGA is orders of magnitude faster than existing state-of-the-art methods.

This work can be extended in several ways. First, we have focused on the SAGA method as the basic iteration loop, but this approach can likely be extended to other proximal incremental schemes such as SGD or ProxSVRG. Second, as mentioned in §3.3, it is an open question whether it is possible to obtain convergence guarantees without any sparsity assumption, as was done for ASAGA.

Acknowledgements

The authors would like to thank our colleagues Damien Garreau, Robert Gower, Thomas Kerdreux, Geoffrey Negiar and Konstantin Mishchenko for their feedback on this manuscript, and Jean-Baptiste Alayrac for support managing the computational resources.

This work was partially supported by a Google Research Award. FP acknowledges support from the chaire *Économie des nouvelles données* with the *data science* joint research initiative with the *fonds AXA pour la recherche*.

References

- Bauschke, Heinz and Combettes, Patrick L. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011.
- Beck, Amir and Teboulle, Marc. *Gradient-based algorithms with applications to signal recovery. Convex Optimization in Signal Processing and Communications*, 2009.
- Davis, Damek, Edmunds, Brent, and Udell, Madeleine. *The sound of APALM clapping: faster nonsmooth nonconvex optimization with stochastic asynchronous PALM*. In *Advances in Neural Information Processing Systems* 29, 2016.
- Defazio, Aaron, Bach, Francis, and Lacoste-Julien, Simon. *SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*. In *Advances in Neural Information Processing Systems*, 2014.
- Gu, Bin, Huo, Zhouyuan, and Huang, Heng. *Asynchronous stochastic block coordinate descent with variance reduction*. *arXiv preprint arXiv:1610.09447v3*, 2016.
- Hsieh, Cho-Jui, Yu, Hsiang-Fu, and Dhillon, Inderjit S. *PASSCoDe: parallel asynchronous stochastic dual coordinate descent*. In *ICML*, 2015.
- Johnson, Rie and Zhang, Tong. *Accelerating stochastic gradient descent using predictive variance reduction*. In *Advances in Neural Information Processing Systems*, 2013.
- Juan, Yuchin, Zhuang, Yong, Chin, Wei-Sheng, and Lin, Chih-Jen. *Field-aware factorization machines for CTR prediction*. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 2016.

- Le Roux, Nicolas, Schmidt, Mark, and Bach, Francis R. [A stochastic gradient method with an exponential convergence rate for finite training sets](#). In *Advances in Neural Information Processing Systems*, 2012.
- Leblond, Rémi, Pedregosa, Fabian, and Lacoste-Julien, Simon. [ASAGA: asynchronous parallel SAGA](#). *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS 2017)*, 2017.
- Liu, Ji and Wright, Stephen J. [Asynchronous stochastic coordinate descent: Parallelism and convergence properties](#). *SIAM Journal on Optimization*, 2015.
- Mania, Horia, Pan, Xinghao, Papailiopoulos, Dimitris, Recht, Benjamin, Ramchandran, Kannan, and Jordan, Michael I. [Perturbed iterate analysis for asynchronous stochastic optimization](#). *SIAM Journal on Optimization*, 2017.
- Meng, Qi, Chen, Wei, Yu, Jingcheng, Wang, Taifeng, Ma, Zhi-Ming, and Liu, Tie-Yan. [Asynchronous stochastic proximal optimization algorithms with variance reduction](#). In *AAAI*, 2017.
- Nesterov, Yurii. *Introductory lectures on convex optimization*. Springer Science & Business Media, 2004.
- Nesterov, Yurii. [Gradient methods for minimizing composite functions](#). *Mathematical Programming*, 2013.
- Niu, Feng, Recht, Benjamin, Re, Christopher, and Wright, Stephen. [Hogwild: A lock-free approach to parallelizing stochastic gradient descent](#). In *Advances in Neural Information Processing Systems*, 2011.
- Peng, Zhimin, Xu, Yangyang, Yan, Ming, and Yin, Wotao. [ARock: an algorithmic framework for asynchronous parallel coordinate updates](#). *SIAM Journal on Scientific Computing*, 2016.
- Reddi, Sashank J, Hefny, Ahmed, Sra, Suvrit, Póczos, Barnabas, and Smola, Alexander J. [On variance reduction in stochastic gradient descent and its asynchronous variants](#). In *Advances in Neural Information Processing Systems*, 2015.
- Schmidt, Mark, Le Roux, Nicolas, and Bach, Francis. [Minimizing finite sums with the stochastic average gradient](#). *Mathematical Programming*, 2016.
- Shalev-Shwartz, Shai and Zhang, Tong. [Stochastic dual coordinate ascent methods for regularized loss minimization](#). *Journal of Machine Learning Research*, 2013.
- Shalev-Shwartz, Shai et al. [Proximal stochastic dual coordinate ascent](#). *arXiv preprint arXiv:1211.2717*, 2012.
- Xiao, Lin and Zhang, Tong. [A proximal stochastic gradient method with progressive variance reduction](#). *SIAM Journal on Optimization*, 2014.
- You, Yang, Lian, Xiangru, Liu, Ji, Yu, Hsiang-Fu, Dhillon, Inderjit S, Demmel, James, and Hsieh, Cho-Jui. [Asynchronous parallel greedy coordinate descent](#). In *Advances In Neural Information Processing Systems*, 2016.
- Yu, Hsiang-Fu, Lo, Hung-Yi, Hsieh, Hsun-Ping, Lou, Jing-Kai, McKenzie, Todd G, Chou, Jung-Wei, Chung, Po-Han, Ho, Chia-Hua, Chang, Chun-Fu, Wei, Yin-Hsuan, et al. [Feature engineering and classifier ensemble for KDD cup 2010](#). In *KDD Cup*, 2010.
- Zhao, Tuo, Yu, Mo, Wang, Yiming, Arora, Raman, and Liu, Han. [Accelerated mini-batch randomized block coordinate descent method](#). In *Advances in neural information processing systems*, 2014.

Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization

Supplementary material

Notations. Throughout the supplementary material we use the following extra notation. We denote by $\langle \cdot, \cdot \rangle_{(i)}$ (resp. $\| \cdot \|_{(i)}$) the scalar product (resp. norm) restricted to blocks in T_i , i.e., $\langle \mathbf{x}, \mathbf{y} \rangle_{(i)} := \sum_{B \in T_i} \langle [\mathbf{x}]_B, [\mathbf{y}]_B \rangle$ and $\|\mathbf{x}\|_{(i)} := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_{(i)}}$. We will also use the following definitions: $\varphi := \sum_{B \in \mathcal{B}} d_B h_B(\mathbf{x})$ and \mathbf{D} is the diagonal matrix defined block-wise as $[\mathbf{D}]_{B,B} = d_B \mathbf{I}_{|B|}$.

The **Bregman divergence** associated with a convex function f for points \mathbf{x}, \mathbf{y} in its domain is defined as:

$$B_f(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle. \quad (7)$$

Note that this is always positive due to the convexity of f .

Appendix A Basic properties

Lemma 1. For any μ -strongly convex function f we have the following inequality:

$$\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 + B_f(\mathbf{x}, \mathbf{y}). \quad (8)$$

Proof. By strong convexity, f verifies the inequality:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{y}), \mathbf{y} - \mathbf{x} \rangle - \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2, \quad (9)$$

for any \mathbf{x}, \mathbf{y} in the domain (see e.g. (Nesterov, 2004)). We then have the equivalences:

$$\begin{aligned} f(\mathbf{x}) &\leq f(\mathbf{y}) + \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle - \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 \\ \iff \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 + f(\mathbf{x}) - f(\mathbf{y}) &\leq \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle \\ \iff \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \underbrace{f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle}_{B_f(\mathbf{x}, \mathbf{y})} &\leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle, \end{aligned} \quad (10)$$

where in the last line we have subtracted $\langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$ from both sides of the inequality. \square

Lemma 2. Let the f_i be L -smooth and convex functions. Then it is verified that:

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2 \leq 2LB_f(\mathbf{x}, \mathbf{y}). \quad (11)$$

Proof. Since each f_i is L -smooth, it is verified (see e.g. Nesterov (2004, Theorem 2.1.5)) that

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2 \leq 2L(f_i(\mathbf{x}) - f_i(\mathbf{y}) - \langle \nabla f_i(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle). \quad (12)$$

The result is obtained by averaging over i . \square

Lemma 3 (Characterization of the proximal operator). Let h be convex lower semicontinuous. Then we have the following characterization of the proximal operator:

$$\mathbf{z} = \text{prox}_{\gamma h}(\mathbf{x}) \iff \frac{1}{\gamma}(\mathbf{x} - \mathbf{z}) \in \partial h(\mathbf{z}). \quad (13)$$

Proof. This is a direct consequence of the first order optimality conditions on the definition of proximal operator, see e.g. (Beck & Teboulle, 2009; Nesterov, 2013). \square

Lemma 4 (Firm non-expansiveness). *Let $\mathbf{x}, \tilde{\mathbf{x}}$ be two arbitrary elements in the domain of φ_i and $\mathbf{z}, \tilde{\mathbf{z}}$ be defined as $\mathbf{z} := \text{prox}_{\varphi_i}(\mathbf{x})$, $\tilde{\mathbf{z}} := \text{prox}_{\varphi_i}(\tilde{\mathbf{x}})$. Then it is verified that:*

$$\langle \mathbf{z} - \tilde{\mathbf{z}}, \mathbf{x} - \tilde{\mathbf{x}} \rangle_{(i)} \geq \|\mathbf{z} - \tilde{\mathbf{z}}\|_{(i)}^2. \quad (14)$$

Proof. By the block-separability of φ_i , the proximal operator is the concatenation of the proximal operators of the blocks. In other words, for any block $B \in T_i$ we have:

$$[\mathbf{z}]_B = \mathbf{prox}_{\gamma\varphi_B}([\mathbf{x}]_B), \quad [\tilde{\mathbf{z}}]_B = \mathbf{prox}_{\gamma\varphi_B}([\tilde{\mathbf{x}}]_B), \quad (15)$$

where φ_B is the restriction of φ_i to B . By firm non-expansiveness of the proximal operator (see e.g. Bauschke & Combettes (2011, Proposition 4.2)) we have that:

$$\langle [\mathbf{z}]_B - [\tilde{\mathbf{z}}]_B, [\mathbf{x}]_B - [\tilde{\mathbf{x}}]_B \rangle \geq \|[\mathbf{z}]_B - [\tilde{\mathbf{z}}]_B\|^2.$$

Summing over the blocks in T_i yields the desired result. \square

Appendix B Sparse Proximal SAGA

This Appendix contains all proofs for Section 2. The main result of this section is Theorem 1, whose proof is structured as follows:

- We start by proving four auxiliary results that will be used later on in the proofs of both synchronous and asynchronous variants. The first is the unbiasedness of key quantities used in the algorithm. The second is a characterization of the solutions of (OPT) in terms of f and φ (defined below) in Lemma 6. The third is a key inequality in Lemma 7 that relates the gradient mapping to other terms that arise in the optimization. The fourth is an upper bound on the variance terms of the gradient estimator, relating it to the Bregman divergence of f and the past gradient estimator terms.
- In Lemma 9, we define an upper bound on the iterates $\|x_t - x^*\|^2$, called a Lyapunov function, and prove an inequality that relates this Lyapunov function value at the current iterate with its value at the previous iterate.
- Finally, in the proof of Theorem 1 we use the previous inequality in terms of the Lyapunov function to prove a geometric convergence of the iterates.

We start by proving the following unbiasedness result, mentioned in §2.

Lemma 5. *Let D_i and φ_i be defined as in §2. Then it is verified that $\mathbf{E}D_i = I_p$ and $\mathbf{E}\varphi_i = h$.*

Proof. Let $B \in \mathcal{B}$ an arbitrary block. We have the following sequence of equalities:

$$\mathbf{E}[D_i]_{B,B} = \frac{1}{n} \sum_{i=1}^n [D_i]_{B,B} = \frac{1}{n} \sum_{i=1}^n d_B \mathbb{1}\{B \in T_i\} I_{|B|} \quad (16)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{n}{n_B} \mathbb{1}\{B \in T_i\} I_{|B|} \quad (17)$$

$$= \left(\frac{1}{n_B} \sum_{i=1}^n \mathbb{1}\{B \in T_i\} \right) I_{|B|} = I_{|B|}, \quad (18)$$

where the last equality comes from the definition of n_B . $\mathbf{E}D_i = I_p$ then follows from the arbitrariness of B .

Similarly, for φ_i we have:

$$\mathbf{E}\varphi_i([x]_B) = \frac{1}{n} \sum_{i=1}^n d_B \mathbb{1}\{B \in T_i\} h_B([x]_B) \quad (19)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{n}{n_B} \mathbb{1}\{B \in T_i\} h_B([x]_B) \quad (20)$$

$$= \left(\frac{1}{n_B} \sum_{i=1}^n \mathbb{1}\{B \in T_i\} \right) h_B([x]_B) = h_B([x]_B), \quad (21)$$

Finally, the result $\mathbf{E}\varphi_i = h$ comes from adding over all blocks. \square

Lemma 6. *x^* is a solution to (OPT) if and only if the following condition is verified:*

$$x^* = \text{prox}_{\gamma\varphi}(x^* - \gamma D \nabla f(x^*)). \quad (22)$$

Proof. By the first order optimality conditions, the solutions to (OPT) are characterized by the sub-differential inclusion $-\nabla f(x^*) \in \partial h(x^*)$. We can then write the following sequence of equiva-

lences:

$$\begin{aligned}
-\nabla f(\mathbf{x}^*) \in \partial h(\mathbf{x}^*) &\iff -D\nabla f(\mathbf{x}^*) \in D\partial h(\mathbf{x}^*) \\
&\quad (\text{multiplying by } D, \text{ equivalence since diagonals are nonzero}) \\
&\iff -D\nabla f(\mathbf{x}^*) \in \partial\varphi(\mathbf{x}^*) \\
&\quad (\text{by definition of } \varphi) \\
&\iff \frac{1}{\gamma}(\mathbf{x}^* - \gamma D\nabla f(\mathbf{x}^*) - \mathbf{x}^*) \in \partial\varphi(\mathbf{x}^*) \\
&\quad (\text{adding and subtracting } \mathbf{x}^*) \\
&\iff \mathbf{x}^* = \mathbf{prox}_{\gamma\varphi}(\mathbf{x}^* - \gamma D\nabla f(\mathbf{x}^*)). \tag{23} \\
&\quad (\text{by Lemma 3})
\end{aligned}$$

Since all steps are equivalences, we have the desired result. \square

The following lemma will be key in the proof of convergence for both the sequential and the parallel versions of the algorithm. With this result, we will be able to bound the product between the gradient mapping and the iterate suboptimality by:

- First, the negative norm of the gradient mapping, which will be key in the parallel setting to cancel out the terms arising from the asynchrony.
- Second, variance terms in $\|\mathbf{v}_i - D_i \nabla f(\mathbf{x}^*)\|^2$ that we will be able to bound by the Bregman divergence using Lemma 2.
- Third and last, a product with terms in $\langle \mathbf{v}_i - D_i \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle$, which taken in expectation gives $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle$ and will allow us to apply Lemma 1 to obtain the contraction terms needed to obtain a geometric rate of convergence.

Lemma 7 (Gradient mapping inequality). *Let \mathbf{x} be an arbitrary vector, \mathbf{x}^* a solution to (OPT), \mathbf{v}_i as defined in (SPS) and $\mathbf{g} = \mathbf{g}(\mathbf{x}, \mathbf{v}_i, i)$ the gradient mapping defined in (3). Then the following inequality is verified for any $\beta > 0$:*

$$\langle \mathbf{g}, \mathbf{x} - \mathbf{x}^* \rangle \geq -\frac{\gamma}{2}(\beta - 2)\|\mathbf{g}\|^2 - \frac{\gamma}{2\beta}\|\mathbf{v}_i - D_i \nabla f(\mathbf{x}^*)\|^2 + \langle \mathbf{v}_i - D_i \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle. \tag{24}$$

Proof. By firm non-expansiveness of the proximal operator (Lemma 4) applied to $\mathbf{z} = \mathbf{prox}_{\gamma\varphi_i}(\mathbf{x} - \gamma\mathbf{v}_i)$ and $\tilde{\mathbf{z}} = \mathbf{prox}_{\gamma\varphi_i}(\mathbf{x}^* - \gamma D\nabla f(\mathbf{x}^*))$ we have:

$$\|\mathbf{z} - \tilde{\mathbf{z}}\|_{(i)}^2 - \langle \mathbf{z} - \tilde{\mathbf{z}}, \mathbf{x} - \gamma\mathbf{v}_i - \mathbf{x}^* + \gamma D\nabla f(\mathbf{x}^*) \rangle_{(i)} \leq 0. \tag{25}$$

By the (SPS) iteration we have $\mathbf{x}^+ = \mathbf{z}$ and by Lemma 3 we have that $[\mathbf{z}]_{T_i} = [\mathbf{x}^*]_{T_i}$, hence the above can be rewritten as

$$\|\mathbf{x}^+ - \mathbf{x}^*\|_{(i)}^2 - \langle \mathbf{x}^+ - \mathbf{x}^*, \mathbf{x} - \gamma\mathbf{v}_i - \mathbf{x}^* + \gamma D\nabla f(\mathbf{x}^*) \rangle_{(i)} \leq 0. \tag{26}$$

We can now write the following sequence of inequalities

$$\begin{aligned}
\langle \gamma \mathbf{g}, \mathbf{x} - \mathbf{x}^* \rangle &= \langle \mathbf{x} - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^* \rangle_{(i)} \quad (\text{by definition and sparsity of } \mathbf{g}) \\
&= \langle \mathbf{x} - \mathbf{x}^+ + \mathbf{x}^* - \mathbf{x}^*, \mathbf{x} - \mathbf{x}^* \rangle_{(i)} \\
&= \|\mathbf{x} - \mathbf{x}^*\|_{(i)}^2 - \langle \mathbf{x}^+ - \mathbf{x}^*, \mathbf{x} - \mathbf{x}^* \rangle_{(i)} \\
&\geq \|\mathbf{x} - \mathbf{x}^*\|_{(i)}^2 - \langle \mathbf{x}^+ - \mathbf{x}^*, 2\mathbf{x} - \gamma \mathbf{v}_i - 2\mathbf{x}^* + \gamma \mathbf{D} \nabla f(\mathbf{x}^*) \rangle_{(i)} + \|\mathbf{x}^+ - \mathbf{x}^*\|_{(i)}^2 \quad (27) \\
&\quad (\text{adding Eq. (26)}) \\
&= \|\mathbf{x} - \mathbf{x}^+\|_{(i)}^2 + \langle \mathbf{x}^+ - \mathbf{x}^*, \gamma \mathbf{v}_i - \gamma \mathbf{D} \nabla f(\mathbf{x}^*) \rangle_{(i)} \quad (\text{completing the square}) \\
&= \|\mathbf{x} - \mathbf{x}^+\|_{(i)}^2 + \langle \mathbf{x} - \mathbf{x}^*, \gamma \mathbf{v}_i - \gamma \mathbf{D} \nabla f(\mathbf{x}^*) \rangle_{(i)} - \langle \mathbf{x} - \mathbf{x}^+, \gamma \mathbf{v}_i - \gamma \mathbf{D} \nabla f(\mathbf{x}^*) \rangle_{(i)} \\
&\quad (\text{adding and subtracting } \mathbf{x}) \\
&\geq \left(1 - \frac{\beta}{2}\right) \|\mathbf{x} - \mathbf{x}^+\|_{(i)}^2 - \frac{\gamma^2}{2\beta} \|\mathbf{v}_i - \mathbf{D} \nabla f(\mathbf{x}^*)\|_{(i)}^2 + \gamma \langle \mathbf{v}_i - \mathbf{D} \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle_{(i)} \\
&\quad (\text{Young's inequality } 2\langle a, b \rangle \leq \frac{\|a\|^2}{\beta} + \beta \|b\|^2, \text{ valid for arbitrary } \beta > 0) \\
&\geq \left(1 - \frac{\beta}{2}\right) \|\mathbf{x} - \mathbf{x}^+\|_{(i)}^2 - \frac{\gamma^2}{2\beta} \|\mathbf{v}_i - \mathbf{D}_i \nabla f(\mathbf{x}^*)\|^2 + \gamma \langle \mathbf{v}_i - \mathbf{D}_i \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \\
&\quad (\text{by definition of } \mathbf{D}_i \text{ and using the fact that } \mathbf{v}_i \text{ is } T_i\text{-sparse}) \\
&= \left(1 - \frac{\beta}{2}\right) \|\gamma \mathbf{g}\|^2 - \frac{\gamma^2}{2\beta} \|\mathbf{v}_i - \mathbf{D}_i \nabla f(\mathbf{x}^*)\|^2 + \gamma \langle \mathbf{v}_i - \mathbf{D} \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle, \quad (28)
\end{aligned}$$

where in the last inequality we have used the fact that \mathbf{g} is T_i -sparse. Finally, dividing by γ both sides yields the desired result. \square

Lemma 8 (Upper bound on the gradient estimator variance). *For arbitrary vectors \mathbf{x} , $(\alpha_i)_{i=0}^n$, and \mathbf{v}_i as defined in (SPS) we have:*

$$\mathbf{E} \|\mathbf{v}_i - \mathbf{D}_i \nabla f(\mathbf{x}^*)\|^2 \leq 4LB_f(\mathbf{x}, \mathbf{x}^*) + 2\mathbf{E} \|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2. \quad (29)$$

Proof. We will now bound the variance terms. For this we have:

$$\begin{aligned}
\mathbf{E} \|\mathbf{v}_i - \mathbf{D} \nabla f(\mathbf{x}^*)\|_{(i)}^2 &= \mathbf{E} \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*) + \nabla f_i(\mathbf{x}^*) - \alpha_i + \mathbf{D}_i \bar{\alpha} - \mathbf{D} \nabla f(\mathbf{x}^*)\|_{(i)}^2 \\
&\leq 2\mathbf{E} \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 + 2\mathbf{E} \|\nabla f_i(\mathbf{x}^*) - \alpha_i - (\mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha})\|_{(i)}^2 \\
&\quad (\text{by inequality } \|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2) \\
&= 2\mathbf{E} \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 + 2\mathbf{E} \|\nabla f_i(\mathbf{x}^*) - \alpha_i\|^2 \\
&\quad - 4\mathbf{E} \langle \nabla f_i(\mathbf{x}^*) - \alpha_i, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle_{(i)} + 2\mathbf{E} \|\mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha}\|_{(i)}^2. \quad (30) \\
&\quad (\text{developing the square})
\end{aligned}$$

We will now simplify the last two terms in the above expression. For the first of the two last terms we have:

$$-4\mathbf{E} \langle \nabla f_i(\mathbf{x}^*) - \alpha_i, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle_{(i)} = -4\mathbf{E} \langle \nabla f_i(\mathbf{x}^*) - \alpha_i, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle \quad (31)$$

$$\begin{aligned}
&\quad (\text{support of first term}) \\
&= -4\langle \nabla f(\mathbf{x}^*) - \bar{\alpha}, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle \\
&= -4\|\nabla f(\mathbf{x}^*) - \bar{\alpha}\|_{\mathbf{D}}^2. \quad (32)
\end{aligned}$$

Similarly, for the last term we have:

$$\begin{aligned}
2\mathbf{E} \|\mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha}\|_{(i)}^2 &= 2\mathbf{E} \langle \mathbf{D}_i \nabla f(\mathbf{x}^*) - \mathbf{D}_i \bar{\alpha}, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle \\
&= 2\langle \nabla f(\mathbf{x}^*) - \bar{\alpha}, \mathbf{D} \nabla f(\mathbf{x}^*) - \mathbf{D} \bar{\alpha} \rangle \\
&\quad (\text{using Lemma 5}) \\
&= 2\|\nabla f(\mathbf{x}^*) - \bar{\alpha}\|_{\mathbf{D}}^2. \quad (33)
\end{aligned}$$

and so the addition of these terms is negative and can be dropped. In all, for the variance terms we have

$$\begin{aligned}\mathbf{E}\|v_i - D\nabla f(\mathbf{x}^*)\|_{(i)}^2 &\leq 2\mathbf{E}\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 + 2\mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2 \\ &\leq 4LB_f(\mathbf{x}, \mathbf{x}^*) + 2\mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2. \quad (\text{by Lemma 2})\end{aligned}\quad (34)$$

□

We now define an upper bound on the quantity that we would like to bound, often called a Lyapunov function, and establish a recursive inequality on this Lyapunov function.

Lemma 9 (Lyapunov inequality). *Let \mathcal{L} be the following c -parametrized function:*

$$\mathcal{L}(\mathbf{x}, \alpha) := \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{c}{n} \sum_{i=1}^n \|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2. \quad (35)$$

Let \mathbf{x}^+ and α^+ be obtained from the Sparse Proximal SAGA updates (SPS). Then we have:

$$\begin{aligned}\mathbf{E}\mathcal{L}(\mathbf{x}^+, \alpha^+) - \mathcal{L}(\mathbf{x}, \alpha) &\leq -\gamma\mu\|\mathbf{x} - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}, \mathbf{x}^*) \\ &\quad + \left(2\gamma^2 - \frac{c}{n}\right) \mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x})\|^2.\end{aligned}\quad (36)$$

Proof. For the first term of \mathcal{L} we have:

$$\begin{aligned}\|\mathbf{x}^+ - \mathbf{x}^*\|^2 &= \|\mathbf{x} - \gamma\mathbf{g} - \mathbf{x}^*\|^2 \quad (\mathbf{g} := \mathbf{g}(\mathbf{x}, v_i, i)) \\ &= \|\mathbf{x} - \mathbf{x}^*\|^2 - 2\gamma\langle \mathbf{g}, \mathbf{x} - \mathbf{x}^* \rangle + \|\gamma\mathbf{g}\|^2 \\ &\leq \|\mathbf{x} - \mathbf{x}^*\|^2 + \gamma^2\|\mathbf{v}_i - D_i\nabla f(\mathbf{x}^*)\|^2 - 2\gamma\langle \mathbf{v}_i - D_i\nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \\ &\quad (\text{by Lemma 7 with } \beta = 1)\end{aligned}\quad (37)$$

Since v_i is an unbiased estimator of the gradient and $\mathbf{E}D_i = I_p$, taking expectations we have:

$$\begin{aligned}\mathbf{E}\|\mathbf{x}^+ - \mathbf{x}^*\|^2 &\leq \|\mathbf{x} - \mathbf{x}^*\|^2 + \gamma^2\mathbf{E}\|\mathbf{v}_i - D_i\nabla f(\mathbf{x}^*)\|^2 - 2\gamma\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \\ &\leq (1 - \gamma\mu)\|\mathbf{x} - \mathbf{x}^*\|^2 + \gamma^2\mathbf{E}\|\mathbf{v}_i - D_i\nabla f(\mathbf{x}^*)\|^2 - 2\gamma B_f(\mathbf{x}, \mathbf{x}^*). \\ &\quad (\text{by Lemma 1})\end{aligned}\quad (38)$$

By using the variance terms bound (Lemma 8) in the previous equation we have:

$$\begin{aligned}\mathbf{E}\|\mathbf{x}^+ - \mathbf{x}^*\|^2 &\leq (1 - \gamma\mu)\|\mathbf{x} - \mathbf{x}^*\|^2 + (4L\gamma^2 - 2\gamma)B_f(\mathbf{x}, \mathbf{x}^*) \\ &\quad + 2\gamma^2\mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2.\end{aligned}\quad (39)$$

We will now bound the second term of the Lyapunov function. We have:

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n \|\alpha_i^+ - \nabla f_i(\mathbf{x}^*)\|^2 &= \left(1 - \frac{1}{n}\right) \mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{1}{n} \mathbf{E}\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)\|^2 \\ &\quad (\text{by definition of } \alpha^+)\end{aligned}\quad (40)$$

$$\begin{aligned}&\leq \left(1 - \frac{1}{n}\right) \mathbf{E}\|\alpha_i - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{2}{n} LB_f(\mathbf{x}, \mathbf{x}^*). \quad (\text{by Lemma 2}) \\ &\quad (41)\end{aligned}$$

Combining Eq. (39) and (40) we have:

$$\begin{aligned}
\mathbf{E}\mathcal{L}(\mathbf{x}^+, \boldsymbol{\alpha}^+) &\leq (1 - \gamma\mu)\|\mathbf{x} - \mathbf{x}^*\|^2 + (4L\gamma^2 - 2\gamma)B_f(\mathbf{x}, \mathbf{x}^*) + 2\gamma^2\mathbf{E}\|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2 \\
&\quad + c \left[\left(1 - \frac{1}{n}\right) \mathbf{E}\|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2 + \frac{1}{n} 2LB_f(\mathbf{x}, \mathbf{x}^*) \right] \\
&= (1 - \gamma\mu)\|\mathbf{x} - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}, \mathbf{x}^*) \\
&\quad + \left(2\gamma^2 - \frac{c}{n}\right) \mathbf{E}\|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2 + c\mathbf{E}\|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2 \\
&= \mathcal{L}(\mathbf{x}, \boldsymbol{\alpha}) - \gamma\mu\|\mathbf{x} - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}, \mathbf{x}^*) \\
&\quad + \left(2\gamma^2 - \frac{c}{n}\right) \mathbf{E}\|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2. \tag{42}
\end{aligned}$$

Finally, subtracting $\mathcal{L}(\mathbf{x}, \boldsymbol{\alpha})$ from both sides yields the desired result. \square

Theorem 1. Let $\gamma = \frac{a}{5L}$ for any $a \leq 1$ and f be μ -strongly convex. Then Sparse Proximal SAGA converges geometrically in expectation with a rate factor of at least $\rho = \frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$. That is, for \mathbf{x}_t obtained after t updates and \mathbf{x}^* the solution to (OPT), we have the bound:

$$\mathbf{E}\|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq (1 - \rho)^t C_0, \quad \text{with } C_0 := \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{1}{5L^2} \sum_{i=1}^n \|\boldsymbol{\alpha}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2.$$

Proof. Let $\bar{H} := \frac{1}{n} \sum_i \|\boldsymbol{\alpha}_i - \nabla f_i(\mathbf{x}^*)\|^2$. By the Lyapunov inequality from Lemma 9, we have:

$$\begin{aligned}
\mathbf{E}\mathcal{L}_{t+1} - (1 - \rho)\mathcal{L}_t &\leq \rho\mathcal{L}_t - \gamma\mu\|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}_t, \mathbf{x}^*) + \left(2\gamma^2 - \frac{c}{n}\right) \bar{H} \\
&= (\rho - \gamma\mu)\|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}_t, \mathbf{x}^*) + \left[2\gamma^2 + c\left(\rho - \frac{1}{n}\right)\right] \bar{H} \\
&\quad \text{(by definition of } \mathcal{L}_t) \\
&\leq (\rho - \gamma\mu)\|\mathbf{x}_t - \mathbf{x}^*\|^2 + \left(4L\gamma^2 - 2\gamma + 2L\frac{c}{n}\right) B_f(\mathbf{x}_t, \mathbf{x}^*) + \left(2\gamma^2 - \frac{2c}{3n}\right) \bar{H} \\
&\quad \text{(choosing } \rho \leq \frac{1}{3n}) \\
&= (\rho - \gamma\mu)\|\mathbf{x}_t - \mathbf{x}^*\|^2 + (10L\gamma^2 - 2\gamma) B_f(\mathbf{x}_t, \mathbf{x}^*) \\
&\quad \text{(choosing } \frac{c}{n} = 3\gamma^2) \\
&\leq \left(\rho - \frac{a\mu}{5L}\right)\|\mathbf{x}_t - \mathbf{x}^*\|^2 \quad \text{(for all } \gamma = \frac{a}{5L}, a \leq 1) \\
&\leq 0. \quad \text{(for } \rho \leq \frac{a}{5} \cdot \frac{\mu}{L}) \tag{43}
\end{aligned}$$

And so we have the bound:

$$\mathbf{E}\mathcal{L}_{t+1} \leq \left(1 - \min\left\{\frac{1}{3n}, \frac{a}{5} \cdot \frac{1}{\kappa}\right\}\right) \mathcal{L}_t \leq \left(1 - \frac{1}{5} \min\left\{\frac{1}{n}, a \cdot \frac{1}{\kappa}\right\}\right) \mathcal{L}_t, \tag{44}$$

where in the last inequality we have used the trivial bound $\frac{1}{3n} \leq \frac{1}{5n}$ merely for clarity of exposition. Chaining expectations from t to 0 we have:

$$\begin{aligned}
\mathbf{E}\mathcal{L}_{t+1} &\leq \left(1 - \frac{1}{5} \min\left\{\frac{1}{n}, a \cdot \frac{1}{\kappa}\right\}\right)^{t+1} \mathcal{L}_0 \\
&= \left(1 - \frac{1}{5} \min\left\{\frac{1}{n}, a \cdot \frac{1}{\kappa}\right\}\right)^{t+1} \left(\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{3a^2}{5^2 L^2} \sum_{i=1}^n \|\boldsymbol{\alpha}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2\right) \\
&\leq \left(1 - \frac{1}{5} \min\left\{\frac{1}{n}, a \cdot \frac{1}{\kappa}\right\}\right)^{t+1} \left(\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{1}{5L^2} \sum_{i=1}^n \|\boldsymbol{\alpha}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2\right) \tag{45} \\
&\quad \text{(since } a \leq 1 \text{ and } 3/5 \leq 1).
\end{aligned}$$

The fact that \mathcal{L}_t is a majorizer of $\|\mathbf{x}_t - \mathbf{x}^*\|^2$ completes the proof. \square

Appendix C ProxASAGA

In this Appendix we provide the proofs for results from Section 3, that is Theorem 2 (the convergence theorem for PROXASAGA) and Corollary 1 (its speedup result).

Notation. Through this section, we use the following shorthand for the gradient mapping: $\mathbf{g}_t := \mathbf{g}(\hat{\mathbf{x}}_t, \hat{\mathbf{v}}_{i_t}^t, i_t)$.

Appendix C.1 Proof outline.

As in the smooth case ($h = 0$), we start by using the definition of \mathbf{x}_{t+1} in Eq. (4) to relate the distance to the optimum in terms of its previous iterates:

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\gamma\langle\hat{\mathbf{x}}_t - \mathbf{x}_t, \mathbf{g}_t\rangle + \gamma^2\|\mathbf{g}_t\|^2 - 2\gamma\langle\hat{\mathbf{x}}_t - \mathbf{x}^*, \mathbf{g}_t\rangle. \quad (46)$$

However, in this case \mathbf{g}_t is not a gradient estimator but a gradient mapping, so we cannot continue as is customary – by using the unbiasedness of the gradient in the $\langle\hat{\mathbf{x}}_t - \mathbf{x}^*, \mathbf{g}_t\rangle$ term together with the strong convexity of f (see Leblond et al. (2017, Section 3.5)).

To circumvent this difficulty, we derive a tailored inequality for the gradient mapping (Lemma 7 in Appendix B), which in turn allows us to use the classical unbiasedness and strong convexity arguments to get the following inequality:

$$\begin{aligned} a_{t+1} \leq & (1 - \frac{\gamma\mu}{2})a_t + \gamma^2\mathbb{E}\|\mathbf{g}_t\|^2 - 2\gamma\mathbb{E}B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) + \underbrace{\gamma\mu\mathbb{E}\|\hat{\mathbf{x}}_t - \mathbf{x}\|^2 + 2\gamma\mathbb{E}\langle\mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t\rangle}_{\text{additional asynchrony terms}} \quad (47) \\ & + \underbrace{\gamma^2(\beta - 2)\mathbb{E}\|\mathbf{g}_t\|^2 + \frac{\gamma^2}{\beta}\mathbb{E}\|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t}\nabla f(\mathbf{x}^*)\|^2}_{\text{additional proximal and variance terms}}, \end{aligned}$$

where $a_t := \mathbb{E}\|\mathbf{x}_t - \mathbf{x}^*\|^2$. Note that since f is strongly convex, $B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) \geq \frac{\mu}{2}\|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2$.

In the smooth setting, one first expresses the additional asynchrony terms as linear combinations of past gradient variance terms ($\mathbb{E}\|\mathbf{g}_u\|^2$) $_{0 \leq u \leq t}$. Then one crucially uses the negative Bregman divergence term to control the variance terms. However, in our current setting, we cannot relate the norm of the gradient mapping $\mathbb{E}\|\mathbf{g}_t\|^2$ to the Bregman divergence (from which h is absent). Instead, we use the negative term $\gamma^2(\beta - 1)\mathbb{E}\|\mathbf{g}_t\|^2$ to control all the ($\mathbb{E}\|\mathbf{g}_u\|^2$) $_{0 \leq u \leq t}$ terms that arise from asynchrony.

The rest of the proof consists in:

- i) expressing the additional asynchrony terms as linear combinations of ($\mathbb{E}\|\mathbf{g}_u\|^2$) $_{0 \leq u \leq t}$, following Leblond et al. (2017, Lemma 1);
- ii) expressing the last variance term, $\|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t}\nabla f(\mathbf{x}^*)\|^2$, as a linear combination of past Bregman divergences (Lemma 8 in Appendix B and Lemma 2 from Leblond et al. (2017));
- iii) defining a Lyapunov function, $\mathcal{L}_t := \sum_{u=0}^t (1 - \rho)^{t-u} a_u$, and proving that it is bounded by a contraction given conditions on the maximum step size and delay.

Appendix C.2 Detailed proof

Theorem 2 (Convergence guarantee and rate of PROXASAGA). *Suppose $\tau \leq \frac{1}{10\sqrt{\Delta}}$. For any step size $\gamma = \frac{a}{L}$ with $a \leq \frac{1}{36} \min\{1, \frac{6\kappa}{\tau}\}$, the inconsistent read iterates of Algorithm 1 converge in expectation at a geometric rate factor of at least: $\rho(a) = \frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$, i.e. $\mathbb{E}\|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2 \leq (1 - \rho)^t \tilde{C}_0$, where \tilde{C}_0 is a constant independent of t ($\approx \frac{n\kappa}{a} C_0$ with C_0 as defined in Theorem 1).*

Proof. In order to get an **initial recursive inequality**, we first unroll the (virtual) update:

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_t - \gamma\mathbf{g}_t - \mathbf{x}^*\|^2 = \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \|\gamma\mathbf{g}_t\|^2 - 2\gamma\langle\mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^*\rangle \\ &= \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \|\gamma\mathbf{g}_t\|^2 - 2\gamma\langle\mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}^*\rangle + 2\gamma\langle\mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t\rangle, \end{aligned} \quad (48)$$

and then apply Lemma 7 with $\mathbf{x} = \hat{\mathbf{x}}_t$ and $\mathbf{v} = \hat{\mathbf{v}}_{i_t}^t$. Note that in this case we have $\mathbf{g} = \mathbf{g}_t$ and $\langle \cdot \rangle_{(i)} = \langle \cdot \rangle_{(i_t)}$.

$$\begin{aligned}
\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 &\leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\gamma \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle + \gamma^2 \|\mathbf{g}_t\|^2 + \gamma^2(\beta - 2) \|\mathbf{g}_t\|^2 \\
&\quad + \frac{\gamma^2}{\beta} \|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D} \nabla f(\mathbf{x}^*)\|_{(i_t)}^2 - 2\gamma \langle \hat{\mathbf{v}}_{i_t}^t - \mathbf{D} \nabla f(\mathbf{x}^*), \hat{\mathbf{x}}_t - \mathbf{x}^* \rangle_{(i_t)} \\
&= \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\gamma \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle + \gamma^2(\beta - 1) \|\mathbf{g}_t\|^2 \\
&\quad + \frac{\gamma^2}{\beta} \|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t} \nabla f(\mathbf{x}^*)\|^2 - 2\gamma \langle \hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t} \nabla f(\mathbf{x}^*), \hat{\mathbf{x}}_t - \mathbf{x}^* \rangle. \quad (49) \\
&\quad (\text{as } [\hat{\mathbf{v}}_{i_t}^t]_{T_{i_t}} = \hat{\mathbf{v}}_{i_t}^t)
\end{aligned}$$

We now use the property that i_t is independent of $\hat{\mathbf{x}}_t$ (which we enforce by reading $\hat{\mathbf{x}}_t$ before picking i_t , see Section 3), together with the unbiasedness of the gradient update $\hat{\mathbf{v}}_{i_t}^t$ ($\mathbf{E} \hat{\mathbf{v}}_{i_t}^t = \nabla f(\hat{\mathbf{x}}_t)$) and the definition of \mathbf{D} to simplify the following expression as follows:

$$\begin{aligned}
\mathbf{E} \langle \hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t} \nabla f(\mathbf{x}^*), \hat{\mathbf{x}}_t - \mathbf{x}^* \rangle &= \langle \nabla f(\hat{\mathbf{x}}_t) - \nabla f(\mathbf{x}^*), \hat{\mathbf{x}}_t - \mathbf{x}^* \rangle \\
&\geq \frac{\mu}{2} \|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2 + B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*), \quad (50)
\end{aligned}$$

where the last inequality comes from Lemma 1. Taking conditional expectations on (49) we get:

$$\mathbf{E} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\gamma \mathbf{E} \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle + \gamma^2(\beta - 1) \mathbf{E} \|\mathbf{g}_t\|^2 \quad (51)$$

$$\begin{aligned}
&+ \frac{\gamma^2}{\beta} \mathbf{E} \|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t} \nabla f(\mathbf{x}^*)\|^2 - \gamma \mu \|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2 - 2\gamma B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) \\
&\leq (1 - \frac{\gamma \mu}{2}) \|\mathbf{x}_t - \mathbf{x}^*\|^2 + 2\gamma \mathbf{E} \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle + \gamma^2(\beta - 1) \mathbf{E} \|\mathbf{g}_t\|^2 \\
&\quad + \frac{\gamma^2}{\beta} \mathbf{E} \|\hat{\mathbf{v}}_{i_t}^t - \mathbf{D}_{i_t} \nabla f(\mathbf{x}^*)\|^2 + \gamma \mu \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|^2 - 2\gamma B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) \\
&\quad (\text{using } \|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2 \text{ on } \|\mathbf{x}_t - \hat{\mathbf{x}}_t + \hat{\mathbf{x}}_t - \mathbf{x}^*\|^2) \\
&\leq (1 - \frac{\gamma \mu}{2}) \|\mathbf{x}_t - \mathbf{x}^*\|^2 + \gamma^2(\beta - 1) \mathbf{E} \|\mathbf{g}_t\|^2 + \gamma \mu \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|^2 + 2\gamma \mathbf{E} \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle \\
&\quad - 2\gamma B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) + \frac{4\gamma^2 L}{\beta} B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) + \frac{2\gamma^2}{\beta} \mathbf{E} \|\hat{\alpha}_{i_t}^t - \nabla f_{i_t}(\mathbf{x}^*)\|^2. \quad (52) \\
&\quad (\text{using Lemma 8 on the variance terms})
\end{aligned}$$

Since we also have:

$$\hat{\mathbf{x}}_t - \mathbf{x}_t = \gamma \sum_{u=(t-\tau)_+}^{t-1} \mathbf{G}_u^t \mathbf{g}(\hat{\mathbf{x}}_u, \hat{\alpha}^u, i_u), \quad (53)$$

the effect of asynchrony for the perturbed iterate updates was already derived in a very similar setup in Leblond et al. (2017). We re-use the following bounds from their Appendix C.4:⁷

$$\mathbb{E} \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|^2 \leq \gamma^2(1 + \sqrt{\Delta} \tau) \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E} \|\mathbf{g}_u\|^2, \quad \text{Leblond et al. (2017, Eq. (48))} \quad (54)$$

$$\mathbb{E} \langle \mathbf{g}_t, \hat{\mathbf{x}}_t - \mathbf{x}_t \rangle \leq \frac{\gamma \sqrt{\Delta}}{2} \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E} \|\mathbf{g}_u\|^2 + \frac{\gamma \sqrt{\Delta} \tau}{2} \mathbb{E} \|\mathbf{g}_t\|^2. \quad \text{Leblond et al. (2017, Eq. (46)).} \quad (55)$$

⁷The appearance of the sparsity constant Δ is coming from the crucial property that $\mathbf{E} \|\mathbf{x}\|_{(i)}^2 \leq \Delta \|\mathbf{x}\|^2 \forall \mathbf{x} \in \mathbb{R}^p$ (see Eq. (39) in Leblond et al. (2017), where they use the notation $\|\cdot\|_i$ for our $\|\cdot\|_{(i)}$).

Because the updates on α are the same for PROXASAGA as for ASAGA, we can re-use the same argument arising in the proof of [Leblond et al. \(2017, Lemma 2\)](#) to get the following bound on $\mathbb{E}\|\hat{\alpha}_{i_t}^t - \nabla f_{i_t}(\mathbf{x}^*)\|^2$:

$$\mathbb{E}\|\hat{\alpha}_{i_t}^t - \nabla f_{i_t}(\mathbf{x}^*)\|^2 \leq \underbrace{\frac{2L}{n} \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)+} \mathbb{E}B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*) + 2L\left(1 - \frac{1}{n}\right)^{(t-\tau)+} \tilde{e}_0}_{\text{Henceforth denoted } H_t}, \quad (56)$$

where $\tilde{e}_0 := \frac{1}{2L} \mathbb{E}\|\alpha_i^0 - f'_i(\mathbf{x}^*)\|^2$. This bound is obtained by analyzing which gradient could be the source of $\hat{\alpha}_{i_t}$ in the past (taking in consideration the inconsistent writes), and then applying Lemma 2 on the $\mathbb{E}\|\nabla f(\hat{\mathbf{x}}_u) - \nabla f(\mathbf{x}^*)\|^2$ terms, explaining the presence of $B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*)$ terms.⁸ The inequality (56) corresponds to Eq. (56) and (57) in [Leblond et al. \(2017\)](#).

By taking the full expectation of (52) and plugging the above inequalities back, we obtain an inequality similar to [Leblond et al. \(2017, Master inequality \(28\)\)](#) which describes how the error terms $a_t := \mathbb{E}\|\mathbf{x}_t - \mathbf{x}^*\|^2$ of the virtual iterates are related:

$$\begin{aligned} a_{t+1} &\leq \left(1 - \frac{\gamma\mu}{2}\right)a_t + \frac{4\gamma^2 L}{\beta} \left(1 - \frac{1}{n}\right)^{(t-\tau)+} \tilde{e}_0 \\ &\quad + \gamma^2 \left[\beta - 1 + \sqrt{\Delta}\tau\right] \mathbb{E}\|\mathbf{g}_t\|^2 + \left[\gamma^2\sqrt{\Delta} + \gamma^3\mu(1 + \sqrt{\Delta}\tau)\right] \sum_{u=(t-\tau)+}^t \mathbb{E}\|\mathbf{g}_u\|^2 \\ &\quad - 2\gamma\mathbb{E}B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) + \frac{4\gamma^2 L}{\beta} \mathbb{E}B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) + \frac{4\gamma^2 L}{\beta n} H_t. \end{aligned} \quad (57)$$

We now have a promising inequality with a contractive term and several quantities that we need to bound. In order to achieve our final result, we introduce the same Lyapunov function as in [Leblond et al. \(2017\)](#):

$$\mathcal{L}_t := \sum_{u=0}^t (1 - \rho)^{t-u} a_u,$$

where ρ is a target rate factor for which we will provide a value later on. Proving that this Lyapunov function is bounded by a contraction will finish our proof. We have:

$$\begin{aligned} \mathcal{L}_{t+1} &= \sum_{u=0}^{t+1} (1 - \rho)^{t+1-u} a_u = (1 - \rho)^{t+1} a_0 + \sum_{u=1}^{t+1} (1 - \rho)^{t+1-u} a_u \\ &= (1 - \rho)^{t+1} a_0 + \sum_{u=0}^t (1 - \rho)^{t-u} a_{u+1}. \end{aligned} \quad (58)$$

We now plug our new bound on a_{t+1} , (57):

$$\begin{aligned} \mathcal{L}_{t+1} &\leq (1 - \rho)^{t+1} a_0 + \sum_{u=0}^t (1 - \rho)^{t-u} \left[\left(1 - \frac{\gamma\mu}{2}\right)a_u + \frac{4\gamma^2 L}{\beta} \left(1 - \frac{1}{n}\right)^{(u-\tau)+} \tilde{e}_0 \right. \\ &\quad + \gamma^2 (\beta - 1 + \sqrt{\Delta}\tau) \mathbb{E}\|\mathbf{g}_u\|^2 \\ &\quad + (\gamma^2\sqrt{\Delta} + \gamma^3\mu(1 + \sqrt{\Delta}\tau)) \sum_{v=(u-\tau)+}^u \mathbb{E}\|\mathbf{g}_v\|^2 \\ &\quad \left. - 2\gamma\mathbb{E}B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*) + \frac{4\gamma^2 L}{\beta} \mathbb{E}B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*) + \frac{4\gamma^2 L}{\beta n} H_u \right]. \end{aligned} \quad (59)$$

After regrouping similar terms, we get:

$$\mathcal{L}_{t+1} \leq (1 - \rho)^{t+1} (a_0 + A\tilde{e}_0) + (1 - \frac{\gamma\mu}{2})\mathcal{L}_t + \sum_{u=0}^t s_u^t \mathbb{E}\|\mathbf{g}_u\|^2 + \sum_{u=1}^t r_u^t \mathbb{E}B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*). \quad (60)$$

⁸Note that [Leblond et al. \(2017\)](#) analyzed the unconstrained scenario, and so $B_f(\hat{\mathbf{x}}_u, \mathbf{x}^*)$ is replaced by the simpler $f(\hat{\mathbf{x}}_u) - f(\mathbf{x}^*)$ in their bound.

Now, provided that we can prove that under certain conditions the s_u^t and r_u^t terms are all negative (and that the A term is not too big), we can drop them from the right-hand side of (60) which will allow us to finish the proof.

Let us compute these terms. Let $q := \frac{1-1/n}{1-\rho}$ and we assume in the rest that $\rho < 1/n$.

Computing A . We have:

$$\begin{aligned} \frac{4\gamma^2 L}{\beta} \sum_{u=0}^t (1-\rho)^{t-u} \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} &\leq \frac{4\gamma^2 L}{\beta} (1-\rho)^t (1-\rho)^{-\tau} \left(\tau + 1 + \frac{1}{1-q}\right) \\ &\quad \text{from Leblond et al. (2017, Eq (75))} \\ &= (1-\rho)^{t+1} \underbrace{\frac{4\gamma^2 L}{\beta} (1-\rho)^{-\tau-1} \left(\tau + 1 + \frac{1}{1-q}\right)}_{:=A}. \end{aligned} \quad (61)$$

Computing s_u^t . Since we have:

$$\sum_{u=0}^t (1-\rho)^{t-u} \sum_{v=(u-\tau)_+}^{u-1} \mathbb{E} \|g_u\|^2 \leq \tau (1-\rho)^{-\tau} \sum_{u=0}^t (1-\rho)^{t-u} \mathbb{E} \|g_u\|^2, \quad (62)$$

we have for all $0 \leq u \leq t$:

$$s_u^t \leq (1-\rho)^{t-u} \left[\gamma^2 (\beta - 1 + \sqrt{\Delta} \tau) + \tau (1-\rho)^{-\tau} (\gamma^2 \sqrt{\Delta} + \gamma^3 \mu (1 + \sqrt{\Delta} \tau)) \right]. \quad (63)$$

Computing r_u^t . To analyze these quantities, we need to compute: $\sum_{u=0}^t (1-\rho)^{t-u} \sum_{v=1}^{u-1} (1 - \frac{1}{n})^{(u-2\tau-v-1)_+}$. Fortunately, this is already done in Leblond et al. (2017, Eq (66)), and thus we know that for all $1 \leq u \leq t$:

$$r_u^t \leq (1-\rho)^{t-u} \left[-2\gamma + \frac{4\gamma^2 L}{\beta} + \frac{4L\gamma^2}{n\beta} (1-\rho)^{-2\tau-1} \left(2\tau + \frac{1}{1-q}\right) \right], \quad (64)$$

recalling that $q := \frac{1-1/n}{1-\rho}$ and that we assumed $\rho < \frac{1}{n}$.

We now need some assumptions to further analyze these quantities. We make simple choices for simplicity, though a tighter analysis is possible. To get manageable (and simple) constants, we follow Leblond et al. (2017, Eq. (82) and (83)) and assume:

$$\rho \leq \frac{1}{4n}; \quad \tau \leq \frac{n}{10}. \quad (65)$$

This tells us:

$$\begin{aligned} \frac{1}{1-q} &\leq \frac{4n}{3} \\ (1-\rho)^{-k\tau-1} &\leq \frac{4}{3} \quad \text{for } 0 \leq k \leq 2. \quad (\text{using Bernouilli's inequality}) \end{aligned}$$

Additionally, we set $\beta = \frac{1}{2}$. Equation (63) thus becomes:

$$s_u^t \leq \gamma^2 (1-\rho)^{t-u} \left[-\frac{1}{2} + \sqrt{\Delta} \tau + \frac{4}{3} (\sqrt{\Delta} \tau + \gamma \mu \tau (1 + \sqrt{\Delta} \tau)) \right]. \quad (66)$$

We see that for s_u^t to be negative, we need $\tau = \mathcal{O}(\frac{1}{\sqrt{\Delta}})$. Let us assume that $\tau \leq \frac{1}{10\sqrt{\Delta}}$. We then get:

$$s_u^t \leq \gamma^2 (1-\rho)^{t-u} \left[-\frac{1}{2} + \frac{1}{10} + \frac{4}{30} + \gamma \mu \tau \frac{4}{3} \frac{11}{10} \right]. \quad (67)$$

Thus, the condition under which all s_u^t are negative boils down to:

$$\gamma \mu \tau \leq \frac{2}{11}. \quad (68)$$

Now looking at the r_u^t terms given our assumptions, the inequality (64) becomes:

$$\begin{aligned} r_u^t &\leq (1 - \rho)^{t-u} \left[-2\gamma + 8\gamma^2 L + \frac{8\gamma^2 L}{n} \frac{4}{3} \left(\frac{n}{5} + \frac{4n}{3} \right) \right] \\ &\leq (1 - \rho)^{t-u} (-2\gamma + 36\gamma^2 L). \end{aligned} \quad (69)$$

The condition for all r_u^t to be negative then can be simplified down to:

$$\gamma \leq \frac{1}{18L}. \quad (70)$$

We now have a promising inequality for proving that our Lyapunov function is bounded by a contraction. However we have defined \mathcal{L}_t in terms of the virtual iterate \mathbf{x}_t , which means that our result would only hold for a given T fixed in advance, as is the case in [Mania et al. \(2017\)](#). Fortunately, we can use the same trick as in [Leblond et al. \(2017, Eq. \(97\)\)](#): we simply add $\gamma B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*)$ to both sides in (60). r_t^t is replaced by $r_t^t + \gamma$, which makes for a slightly worse bound on γ to ensure linear convergence:

$$\gamma \leq \frac{1}{36L}. \quad (71)$$

For this small cost, we get a contraction bound on $B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*)$, and thus by the strong convexity of f (see (9)) we get a contraction bound for $\mathbb{E}\|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2$.

Recap. Let us use $\rho = \frac{1}{4n}$ and $\gamma := \frac{a}{L}$. Then the conditions (68) and (71) on the step size γ reduce to:

$$a \leq \frac{1}{36} \min\{1, \frac{72}{11} \frac{\kappa}{\tau}\}. \quad (72)$$

Moreover, the condition:

$$\tau \leq \frac{1}{10\sqrt{\Delta}} \quad (73)$$

is sufficient to also ensure that (65) is satisfied as $\Delta \in [\frac{1}{n}, 1]$, and thus $\frac{1}{\sqrt{\Delta}} \leq \sqrt{n} \leq n$.

Thus under the conditions (72) and (73), we have that all s_u^t and r_u^t terms are negative and we can rewrite the recurrent step of our Lyapunov function as:

$$\mathcal{L}_{t+1} \leq \gamma \mathbb{E} B_f(\hat{\mathbf{x}}_t) + \mathcal{L}_{t+1} \leq (1 - \rho)^{t+1} (a_0 + A\tilde{e}_0) + (1 - \frac{\gamma\mu}{2}) \mathcal{L}_t. \quad (74)$$

By unrolling the recursion (74), we can carefully combine the effect of the geometric term $(1 - \rho)$ with the one of $(1 - \frac{\gamma\mu}{2})$. This was already done in [Leblond et al. \(2017, Apx C.9, Eq. \(101\) to \(103\)\)](#), with a trick to handle various boundary cases, yielding the overall rate:

$$\mathbb{E} B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) \leq (1 - \rho^*)^{t+1} \hat{C}_0, \quad (75)$$

where $\rho^* = \min\{\frac{1}{5n}, a\frac{2}{5\kappa}\}$ (that we simplified to $\rho^* = \frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$ in the theorem statement). To get the final constant, we need to bound A . We have:

$$\begin{aligned} A &= \frac{4\gamma^2 L}{\beta} (1 - \rho)^{-\tau-1} (\tau + 1 + \frac{1}{1-q}) \\ &\leq 8\gamma^2 L \frac{4}{3} (\frac{n}{10} + 1 + \frac{4n}{3}) \\ &\leq 26\gamma^2 L n \\ &\leq \gamma n. \end{aligned} \quad (76)$$

This is the same bound on A that was used by [Leblond et al. \(2017\)](#) and so we obtain the same constant as their Eq. (104):

$$\hat{C}_0 := \frac{21n}{\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \gamma \frac{n}{2L} \mathbb{E} \|\boldsymbol{\alpha}_i^0 - \nabla f_i(\mathbf{x}^*)\|^2). \quad (77)$$

Note that $\hat{C}_0 = \mathcal{O}(\frac{n}{\gamma} C_0)$ with C_0 defined as in Theorem 1.

Now, using the strong convexity of f via (9), we get:

$$\mathbb{E}\|\hat{\mathbf{x}}_t - \mathbf{x}^*\|^2 \leq \frac{2}{\mu} \mathbb{E}B_f(\hat{\mathbf{x}}_t, \mathbf{x}^*) \leq (1 - \rho^*)^{t+1} \tilde{C}_0, \quad (78)$$

where $\tilde{C}_0 = \mathcal{O}(\frac{n\kappa}{a}C_0)$.

This finishes the proof for Theorem 2. \square

Corollary 3 (Speedup). *Suppose $\tau \leq \frac{1}{10\sqrt{\Delta}}$. If $\kappa \geq n$, then using the step size $\gamma = 1/36L$, PROXASAGA converges geometrically with rate factor $\Omega(\frac{1}{\kappa})$. If $\kappa < n$, then using the step size $\gamma = 1/36n\mu$, PROXASAGA converges geometrically with rate factor $\Omega(\frac{1}{n})$. In both cases, the convergence rate is the same as Sparse Proximal SAGA and PROXASAGA is thus linearly faster than its sequential counterpart up to a constant factor. Note that in both cases the step size does not depend on τ .*

Furthermore, if $\tau \leq 6\kappa$, we can use a universal step size of $\Theta(1/L)$ to get a similar rate for PROXASAGA than Sparse Proximal SAGA, thus making it adaptive to local strong convexity since the knowledge of κ is not required.

Proof. If $\kappa \geq n$, the rate factor of Sparse Proximal SAGA is $1/\kappa$. To get the same rate factor, we need to choose $a = \Omega(1)$, which we can fortunately do since $\kappa \geq n \geq \sqrt{n} \geq 10\frac{1}{10\sqrt{\Delta}} \geq 10\tau$.

If $\kappa < n$, then the rate factor of Sparse Proximal SAGA is $1/n$. Any choice of a bigger than $\Omega(\kappa/n)$ gives us the same rate factor for PROXASAGA. Since $\tau \leq \sqrt{n}/10$ we can pick such an a without violating the condition of Theorem 2. \square

Appendix D Comparison with related work

In this section, we relate our theoretical results and proof technique with the related literature.

Speedups. Our speedup regimes are comparable with the best ones obtained in the smooth case, including Niu et al. (2011); Reddi et al. (2015), even though unlike these papers, we support inconsistent reads and nonsmooth objective functions. The one exception is Leblond et al. (2017), where the authors prove that their algorithm, ASAGA, can obtain a linear speedup even without sparsity in the well-conditioned regime. In contrast, PROXASAGA always requires some sparsity. Whether this property for smooth objective functions could be extended to the composite case remains an open problem.

Coordinate Descent. We compare our approach for composite objective functions to its most natural competitor: ASYSPCD (Liu & Wright, 2015), an asynchronous stochastic coordinate descent algorithm. While ASYSPCD also exhibits linear speedups, subject to a condition on τ , one has to be especially careful when trying to compare these conditions.

First, while in theory the iterations of both algorithms have the same cost, in practice various tricks are introduced to save on computation, yielding different costs per updates.⁹ Second, the bound on τ for the coordinate descent algorithm depends on p , the dimensionality of the problem, whereas ours involves n , the number of data points. Third, a more subtle issue is that τ is not affected by the same quantities for both algorithms.¹⁰ See Appendix D.1 for a more detailed explanation of the differences between the bounds.

In the best case scenario (where the components of the gradient are uncorrelated, a somewhat unrealistic setting), ASYSPCD can get a near-linear speedup for τ as big as $\sqrt[4]{p}$. Our result states that $\tau = \mathcal{O}(1/\sqrt{\Delta})$ is necessary for a linear speedup. This means in case $\Delta \leq 1/\sqrt{p}$ our bound is better than the one obtained for ASYSPCD. Recalling that $1/n \leq \Delta \leq 1$, it appears that PROXASAGA is favored when n is bigger than \sqrt{p} whereas ASYSPCD may have a better bound otherwise, though this comparison should be taken with a grain of salt given the assumptions we had to make to arrive at comparable quantities.

Furthermore, one has to note that while Liu & Wright (2015) use the classical labeling scheme inherited from Niu et al. (2011), they still assume in their proof that the i_t are uniformly distributed and that their gradient estimators are conditionally unbiased – though neither property is verified in the general asynchronous setting. Finally, we note that ASYSPCD (as well as its incremental variant Async-PROXSVRCD) assumes that the computation and assignment of the proximal operator is an atomic step, while we do not make such assumption.

SVRG. The Async-ProxSVRG algorithm of Meng et al. (2017) also exhibits theoretical linear speedups subject to the same condition as ours. However, the analyzed algorithm uses dense updates and consistent read and writes. Although they make the analysis easier, these two factors introduce costly bottlenecks and prevent linear speedups in running time. Furthermore, here again the classical labeling scheme is used together with the unverified conditional unbiasedness condition.

Doubly stochastic algorithms. The Async-PROXSVRCD algorithm from Meng et al. (2017); Gu et al. (2016) has a maximum allowable stepsize¹¹ that is in $\mathcal{O}(1/pL)$, whereas the maximum step size for PROXASAGA is in $\Omega(1/L)$, so can be up to p times bigger. Consequently, PROXASAGA enjoys much faster theoretical convergence rates. Unfortunately, we could not find a condition for linear speedups to compare to. We also note that their algorithm is not appropriate in a sparse features setting. This is illustrated in an empirical comparison in Appendix F where we see that

⁹For PROXASAGA the relevant quantity becomes the average number of features per data point. For ASYSPCD it is rather the average number of data points per feature. In both cases the tricks involved are not covered by the theory.

¹⁰To make sure τ is the same quantity for both algorithms, we have to assume that the iteration costs are homogeneous.

¹¹To the best of our understanding, noting that extracting an interpretable bound from the given theoretical results was difficult. Furthermore, it appears that the proof technique may still have significant issues: for example, the “fully lock-free” assumption of Gu et al. (2016) allows for overwrites, and is thus incompatible with their framework of analysis, in particular their Eq. (8).

their convergence in number of iterations is orders of magnitude slower than appropriate algorithms like SAGA or PROXASAGA.

Appendix D.1 Comparison of bounds with Liu & Wright (2015)

Iteration costs. For both PROXASAGA and ASYSPCD, the average cost of an iteration is $\mathcal{O}(n\bar{S})$ (where \bar{S} is the average support size). In the case of PROXASAGA (see Algorithm 1), at each iteration the most costly operation is the computation of $\bar{\alpha}$, while in the general case we need to compute a full gradient for ASYSPCD.

In order to reduce these prohibitive computation costs, several tricks are introduced. Although they lead to much improved empirical performance, it should be noted that in both cases these tricks are not covered by the theory. In particular, the unbiasedness condition can be violated.

In the case of PROXASAGA, we store the average gradient term $\bar{\alpha}$ in shared memory. The cost of each iteration then becomes the size of the extended support of the partial gradient selected at random at this iteration, hence it is in $\mathcal{O}(\Delta_l)$, where $\Delta_l := \max_{i=1..n} |T_i|$.

For ASYSPCD, following Peng et al. (2016) we can store intermediary quantities for specific losses (e.g. ℓ_1 -regularized logistic regression). The cost of an iteration then becomes the number of data points whose extended support includes the coordinate selected at random at this iteration, hence it is in $\mathcal{O}(n\Delta)$.

The relative difference in update cost of both algorithms then depends heavily on the data matrix: if the partial gradients usually have a extended support but coordinates belong to few of them (this can be the case if $n \ll p$ for example), then the iterations of ASYSPCD can be cheaper than those of PROXASAGA. Conversely, if data points usually have small extended support but coordinates belong to many of them (which can happen when $p \ll n$ for example), then the updates of PROXASAGA are the cheaper ones.

Dependency of τ on the data matrix. In the case of PROXASAGA the sizes of the extended support of each data point are important – they are directly linked to the cost of each iteration. Identical iteration costs for each data point do not influence τ , whereas heterogeneous costs may cause τ to increase substantially. In contrast, in the case of ASYSPCD, the relevant parts of the data matrix are the number of data points each dimension touches – for much the same reason. In the bipartite graph between data points and dimensions, either the left or the right degrees matter for τ , depending on which algorithm you choose.

In order to compare their respective bounds, we have to make the assumption that the iteration costs are homogeneous, which means that each data point has the same support size and each dimension is active in the same number of data points. This implies that τ is the same quantity for both algorithms.

Best case scenario bound for AsySPCD. The result obtained in Liu & Wright (2015) states that if $\tau^2 \Lambda = \mathcal{O}(\sqrt{p})$, ASYSPCD can get a near-linear speedup (where Λ is a measure of the interactions between the components of the gradient, with $1 \leq \Lambda \leq \sqrt{p}$). In the best possible scenario where $\Lambda = 1$ (which means that the coordinates of the gradients are completely uncorrelated), τ can be as big as $\sqrt[4]{p}$.

Appendix E Implementation details

Initialization. In the Sparse Proximal SAGA algorithm and its asynchronous variant, PROXASAGA, the vector \mathbf{x} can be initialized arbitrarily. The memory terms α_i can be initialized to any vector that verifies $\text{supp}(\alpha_i) = \text{supp}(\nabla f_i)$. In practice we found that the initialization $\alpha_i = \mathbf{0}$ is very fast to set up and often outperforms more costly initializations.

With this initialization, the gradient approximation before the first update of the memory terms becomes $\nabla f_i(\mathbf{x}) + D_i \bar{\alpha}$. Since most of the values in α are zero, $\bar{\alpha}$ will tend to be small compared to $\nabla f_i(\mathbf{x})$, and so the gradient estimate is very close to the SGD estimate $\nabla f_i(\mathbf{x})$. The SGD approximation is known to have a very fast initial convergence (which, in light of Figure 1, our method inherits) and has even been used as a heuristic to use during the first epoch of variance reduced methods (Schmidt et al., 2016).

The initialization of coefficients \mathbf{x}_0 was always set to zero.

Exact regularization. Computing the gradient of a smooth regularization such as the squared ℓ_2 penalty of Eq. (6) is independent of n and so we can use the exact regularizer in the update of the coefficients instead of storing it in α , which would also destroy the compressed storage of the memory terms described below. In practice we use this “exact regularization”, multiplied by D_i to preserve the sparsity pattern.

Assuming a squared ℓ_2 regularization term of the form $\frac{\lambda}{2}$, the gradient estimate in (SPS) becomes (note the extra $\lambda \mathbf{x}$)

$$\mathbf{v}_i = \nabla f_i(\mathbf{x}) - \alpha_i + D_i(\bar{\alpha} + \lambda \mathbf{x}). \quad (79)$$

Storage of memory terms. The storage requirements for this method is in the worst case a table of size $n \times p$. However, as for SAG and SAGA, for linearly parametrized loss functions of the form $f_i(\mathbf{x}) = \ell(\mathbf{a}_i^T \mathbf{x})$, where ℓ is some real-valued function and $(\mathbf{a}_i)_{i=1}^n$ are samples associated with the learning problem, this can be reduced to a table of size n (Schmidt et al., 2016, §4.1). This includes popular linear models such as least squares or logistic regression with ℓ the squared or logistic function, respectively.

The reduce storage comes from the fact that in this case the partial gradients have the structure

$$\nabla f_i(\mathbf{x}) = \mathbf{a}_i \underbrace{\ell'(\mathbf{a}_i^T \mathbf{x})}_{\text{scalar}}. \quad (80)$$

Since \mathbf{a}_i is independent of \mathbf{x} , we only need to store the scalar $\ell'(\mathbf{a}_i^T \mathbf{x})$. This decomposition also explains why ∇f_i inherits the sparsity pattern of \mathbf{a}_i .

Atomic updates. Most modern processors have support for atomic operations with minimal overhead. In our case, we implemented a double-precision atomic type using the C++11 atomic features (`std::atomic<double>`). This type implements atomic operations through the compare and swap semantics.

Empirically, we have found it necessary to implement atomic operations at least in the vector α and $\bar{\alpha}$ to reach arbitrary precision. If non-atomic operations are used, the method converges only to a limited precision (around normalized function suboptimality of 10^{-3}), which might be sufficient for some machine learning applications but which we found not satisfying from an optimization point of view.

AsySPCD. Following (Peng et al., 2016) we keep the vector $(\mathbf{a}_i^T \mathbf{x})_{i=1}^n$ in memory and update it at each iteration using atomic updates.

Hardware and software. All experiments were run on a Dell PowerEdge 920 machine with 4 Intel Xeon E7-4830v2 processors with 10 2.2GHz cores each and 384GB 1600 Mhz RAM. The PROXASAGA and ASYSPCD code was implemented on C++ and binded in Python. The FISTA code is implemented in pure Python using NumPY and SciPY for matrix computations (in this case the bottleneck is in large sparse matrix-vector operations for which efficient BLAS routines were used). Our PROXASAGA implementation can be downloaded from <http://github.com/fabianp/ProxASAGA>.

Appendix F Experiments

All datasets used for the experiments were downloaded from the LibSVM dataset suite.¹²

Appendix F.1 Comparison of ProxASAGA with other sequential methods

We provide a comparison between the Sparse Proximal SAGA and related methods in the sequential case. We compare against two methods: the MRBCD method of Zhao et al. (2014) (which forms the basis of Async-PROXSVRCD) and the vanilla implementation of SAGA (Defazio et al., 2014), which does not have the ability to perform sparse updates. We compare in terms of both passes through the data (epochs) and time. We use the same step size for all methods ($1/3L$). Due to the slow convergence of some methods, we use a smaller dataset than the ones used in §4. Dataset RCV1 has $n = 697,641$, $d = 47,236$ and a density of 0.15, while Covtype is a dense dataset with $n = 581,012$, $d = 54$.

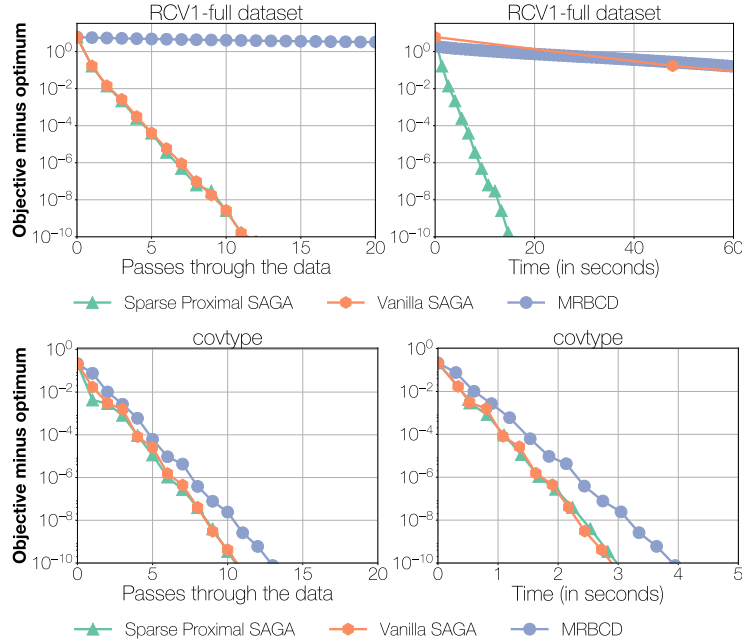


Figure 2: Suboptimality of different sequential algorithms. Each marker represents one pass through the dataset.

We observe that for the convergence behavior in terms of number of passes, Sparse Proximal SAGA performs as well as vanilla SAGA, though the latter requires dense updates at every iteration (Fig. 2 top left). On the other hand, in terms of running time, our implementation of Sparse Proximal SAGA is much more efficient than the other methods for sparse input (Fig. 2 top right). In the case of dense input (Fig. 2 bottom), the three methods perform similarly.

A note on the performance of MRBCD. It may appear surprising that Sparse Proximal SAGA outperforms MRBCD so dramatically on sparse datasets. However, one should note that MRBCD is a doubly stochastic algorithm where both a random data point and a random coordinate are sampled for each iteration. If the data matrix is very sparse, then the probability that the sampled coordinate is in the support of the sampled data point becomes very low. This means that the gradient estimator term only contains the reference gradient term of SVRG, which only changes once per epoch. As a result, this estimator becomes very coarse and produces a slower empirical convergence.

This is reflected in the theoretical results given in Zhao et al. (2014), where the epoch size needed to get linear convergence are k times bigger than the ones required by plain SVRG, where k is the size of the set of blocks of coordinates.

¹²<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Appendix F.2 Theoretical speedups.

In the experimental section, we have shown experimental speedup results where suboptimality was a function of the running time. This measure encompasses both theoretical algorithmic optimization properties and hardware overheads (such as contention of shared memory) which are not taken into account in our analysis.

In order to isolate these two effects, we now plot our speedup results in Figure 3 where suboptimality is a function of the number of iterations; thus, we abstract away any potential hardware overhead. To do so, we implement a global counter which is sparsely updated (every 100 iterations for example) in order not to modify the asynchrony of the system. This counter is used only for plotting purposes and is not needed otherwise. Specifically, we define the theoretical speedup as:

$$\text{theoretical speedup} := (\text{number of cores}) \frac{\text{number of iterations for sequential algorithm}}{\text{total number of iterations for parallel algorithm}}.$$

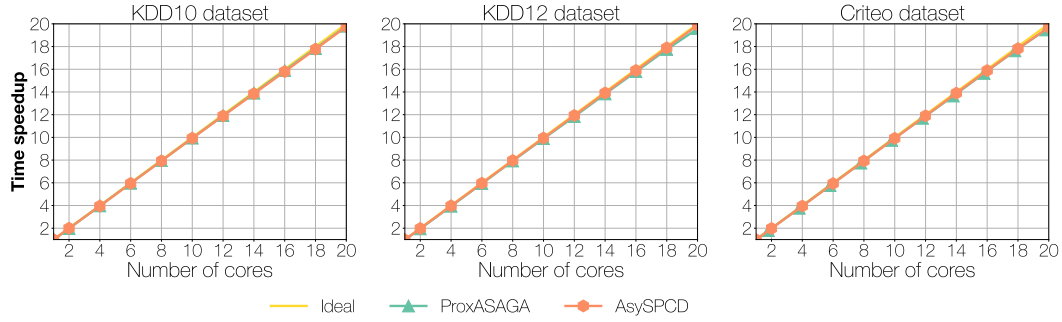


Figure 3: **Theoretical optimization speedups for $\ell_1 + \ell_2$ -regularized logistic regression.** Speedup as measured by the number of iterations required to reach 10^{-5} suboptimality for PROXASAGA and ASYSPCD. In FISTA the iterates are the same with different cores and so matches the “ideal” speedup.

We see clearly that the theoretical speedups obtained by both PROXASAGA and ASYSPCD are linear (i.e. ideal). As we observe worse results in running time, this means that the hardware overheads of asynchronous methods are quite significant.

Appendix F.3 Timing benchmarks

We now provide the time it takes for the different methods with 10 cores to reach a suboptimality of 10^{-10} . All results are in hours.

Dataset	PROXASAGA	ASYSPCD	FISTA
KDD 2010	1.01	13.3	5.2
KDD 2012	0.09	26.6	8.3
Criteo	0.14	33.3	6.6

Appendix F.4 Hyperparameters

The ℓ_1 -regularization parameter λ_2 was chosen as to give around 10% of non-zero features. The exact chosen values are the following: $\lambda_2 = 10^{-11}$ for KDD 2010, $\lambda_2 = 10^{-16}$ for KDD 2012 and $\lambda_2 = 4 \times 10^{-12}$ for Criteo.